

Reed-Muller codes for random erasures and errors

Emmanuel Abbe*

Amir Shpilka[†]Avi Wigderson[‡]

Abstract

This paper studies the parameters for which Reed-Muller (RM) codes over $GF(2)$ can correct random erasures and random errors with high probability, and in particular when can they achieve capacity for these two classical channels. Necessarily, the paper also studies properties of evaluations of multi-variate $GF(2)$ polynomials on random sets of inputs.

For erasures, we prove that RM codes achieve capacity both for very high rate and very low rate regimes. For errors, we prove that RM codes achieve capacity for very low rate regimes, and for very high rates, we show that they can uniquely decode at about square root of the number of errors at capacity.

The proofs of these four results are based on different techniques, which we find interesting in their own right. In particular, we study the following questions about $E(m, r)$, the matrix whose rows are truth tables of all monomials of degree $\leq r$ in m variables. What is the most (resp. least) number of random columns in $E(m, r)$ that define a submatrix having full column rank (resp. full row rank) with high probability? We obtain tight bounds for very small (resp. very large) degrees r , which we use to show that RM codes achieve capacity for erasures in these regimes.

Our decoding from random errors follows from the following novel reduction. For every linear code C of sufficiently high rate we construct a new code C' , also of very high rate, such that for every subset S of coordinates, if C can recover from erasures in S , then C' can recover from errors in S . Specializing this to RM codes and using our results for erasures imply our result on unique decoding of RM codes at high rate.

Finally, two of our capacity achieving results require tight bounds on the weight distribution of RM codes. We obtain such bounds extending the recent [KLP12] bounds from constant degree to linear degree polynomials.

*Program in Applied and Computational Mathematics, and Department of Electrical Engineering, Princeton University, Princeton, USA, eabbe@princeton.edu

[†]Department of Computer Science, Tel-Aviv University, Tel-Aviv, Israel, shpilka@post.tau.ac.il. The research leading to these results has received funding from the European Community's Seventh Framework Programme (FP7/2007-2013) under grant agreement number 257575, and from the Israel Science Foundation (grant number 339/10).

[‡]Institute for Advanced Study, Princeton, USA, avi@ias.edu. This research was partially supported by NSF grant CCF-1412958.

Contents

1	Introduction	1
1.1	Overview	1
1.2	Notation and terminology	4
1.3	Our results	6
1.3.1	Random erasures - the BEC channel	6
1.3.2	Weight distribution and list decoding	7
1.3.3	Random errors - the BSC channel	7
1.4	Proof techniques	8
1.5	Related literature	10
1.6	Organization	11
2	Preliminaries	11
2.1	Basic coding definitions	11
2.2	Equivalent requirements for probabilistic erasures	14
2.3	Basic properties of Reed-Muller codes	16
3	Weight distribution of Reed-Muller codes	17
4	Random submatrices of $E(m, r)$	20
4.1	Random submatrices of $E(m, r)$, for small r , have full column-rank	22
4.2	Random submatrices of $E(m, r)$, for small r , have full row-rank	25
5	Reed-Muller code for erasures	27
5.1	Low-rate regime	27
5.2	High-rate regime	27
6	Reed-Muller code for errors	28
6.1	Low-rate regime	28
6.2	High-rate regime	30
6.2.1	Parity check matrix and parity of patterns	31
6.2.2	The case $r = 1$	32
6.2.3	The degree- r case	34
6.2.4	A general reduction from decoding from errors to decoding from erasures	36
6.2.5	The degree-2 counterexample	37
7	Future directions and open problems	38
A	Proofs of Claim 4.15 and Claim 5.3	42
B	A proof of Lemma 4.10 using hashing	43

1 Introduction

1.1 Overview

We start by giving a high level description of the background and motivation for the problems we study, and of our results.

Reed-Muller (RM) codes were introduced in 1954, first by Muller [Mul54] and shortly after by Reed [Ree54], who also provided a decoding algorithm. They are among the oldest and simplest codes to construct; the codewords are the evaluation vectors of all multivariate polynomials of a given degree bound. More precisely, in an $RM(m, r)$ code over a finite field \mathbb{F} , a message is interpreted as the coefficients of a multivariate polynomial f of degree at most r over \mathbb{F} , and its encoding is simply the vector of evaluations $f(a)$ for all possible assignments $a \in \mathbb{F}^m$ to the variables. Thus, RM codes are linear codes. They have been extensively studied in coding theory, and yet some of their most basic coding-theoretic parameters remain a mystery to date. Specifically, fixing the *rate* of an RM code, while it is easy to compute its tolerance to errors and erasures in the worst-case (or adversarial) model, it has proved extremely difficult to estimate this tolerance for even the simplest models of random errors and erasures. The questions regarding erasures can be interpreted from a learning theory perspective, about interpolating low degree polynomials from lossy or noisy evaluations. The questions regarding errors relate sparse recovery from random Boolean errors. This paper makes some progress on these basic questions.

Reed-Muller codes (over both large and small finite fields) have been extremely influential in the theory of computation, playing a central role in some important developments in several areas. In cryptography, they have been used e.g. in secret sharing schemes [Sha79], instance hiding constructions [BF90] and private information retrieval (see the survey [Gas04]). In the theory of randomness, they have been used in the constructions of many pseudo-random generators and randomness extractors, e.g. [BV10]. These in turn were used for hardness amplification, program testing and eventually in various interactive and probabilistic proof systems, e.g. the celebrated results $NEXP = MIP$ [BFL90], $IP = PSPACE$ [Sha92], $NP = PCP$ [ALM⁺98]. In circuit lower bounds for some low complexity classes one argues that every circuit in the class is close to a codeword, so any function far from the code cannot be computed by such circuits (e.g. [Raz87]). In distributed computing they were used to design fault-tolerant information dispersal algorithms for networks [Rab89]. The hardness of approximation of many optimization problems is greatly improved by the “short code” [BGH⁺12], which uses the optimal testing result of [BKS⁺10]. And the list goes on. Needless to say, the properties used in these works are properties of low-degree polynomials (such interpolation, linearity, partial derivatives, self-reducibility, heredity under various restrictions to variables, etc.), and in some of these cases, specific coding-theoretic perspective such as distance, unique-decoding, list-decoding, local testing and decoding etc. play important roles. Finally, polynomials are basic objects to understand computationally from many perspectives (e.g. testing identities, factoring, learning, etc.), and this study interacts well with the study of coding theoretic questions regarding RM codes.

To discuss the coding-theoretic questions we focus on, and give appropriate perspective, we need some more notation. First, we will restrict attention to binary codes, the most basic case where $\mathbb{F} = \mathbb{F}_2$, the field of two elements¹. To reliably transmit k -bit messages we encode each by an n -bit codeword via a mapping $C : \mathbb{F}_2^k \rightarrow \mathbb{F}_2^n$. We abuse notation and denote by C both the mapping and

¹This seems also the most difficult case for these questions, and we expect our techniques to generalize to larger finite fields.

its image, namely the set of all codewords². The *rate* of C is given by the ratio k/n , capturing the redundancy of the code (the smaller it is, the more redundant it is). A major problem of coding theory is to determine the largest rate for which one can uniquely recover the original message from a *corrupted* codeword (naturally, explicit codes with efficient encoding and decoding algorithms are desirable). This of course depends on the nature of corruption, and we shall deal here with the two most basic ones, erasures (bit-losses) and errors (bit-flips). Curiously, the two seminal papers from the late 1940s giving birth to coding theory, by Shannon [Sha48] and Hamming [Ham50] differ in whether one should consider recovery for *most* corruptions, or from *all* corruptions. In other words, Shannon advocates *average-case* analysis whereas Hamming advocates *worst-case* analysis.

In Hamming’s worst case setting, recovery of the original message must be possible from every corruption of every codeword. In this model there is a single parameter of the code determining recoverability: the *distance* of the code. The *distance* of C is the minimum Hamming distance of any two codewords in C (the *relative distance* is simply the distance normalized by the block-length n). If the distance is d , then we one can uniquely recover from at most d erasures and from $\lfloor (d-1)/2 \rfloor$ errors. This leaves the problem of finding the optimal trade-off between rate and distance, and designing codes which achieve this optimum. While these are still difficult open problems, we know a variety of codes that can simultaneously achieve constant rate and constant relative distance (such codes are often called *asymptotically good*). In contrast, Reed-Muller codes fall far short of that. The rate of $RM(m, r)$ is $\binom{m}{\leq r}/2^m$, while the distance is easily seen to be 2^{m-r} . Thus making any one of them a positive constant makes the other exponentially small in n . In short, from a worst-case perspective, RM codes are pretty bad.

In Shannon’s average-case setting (which we study here), a codeword is subjected to a random corruption, from which recovery should be possible *with high probability*. This random corruption model is called a *channel*, and the best achievable rate is called the *capacity* of the channel. The two most basic ones, the Binary Erasure Channel (BEC) and the Binary Symmetric Channel (BSC), have a parameter p (which may depend on n), and corrupt a message by independently replacing, with probability p , the symbol in each coordinate, with a “lost” symbol in the BEC(p) channel, and with the complementary symbol in the BSC(p) case. Shannon’s original paper already contains the optimal trade-off achievable for these (and many other channels). For *every* p , the capacity of BEC(p) is $1-p$, and the capacity of BSC(p) is $1-h(p)$, where h is the binary entropy function.³ While Shannon shows that random codes achieve this optimal behavior,⁴ explicit and efficiently encodable and decodable codes achieving capacity in both channels⁵ have been obtained [For67], among which are the recent *Polar Codes* [Ari09] that we shall soon discuss.

Do Reed-Muller codes achieve capacity for these natural channels (despite their poor rate-distance trade-off)? The coding theory community seems to believe the answer is positive, and conjectures to that effect were made⁶ in [DG07, Ari08, MHU14]. However, to date, we do not know *any* value of p for which RM codes achieve the capacity for erasures or errors! This paper provides the first progress on this conjecture, resolving it for very low rates and very high rates (namely for polynomials of degrees r which are very small or very large compared to the number of variables m). Our results unfortunately fall short of approaching the cases where the corruption rate p is a

²A code is *linear* if the mapping C is \mathbb{F}_2 -linear, or equivalently if the set of codewords C is a linear subspace of \mathbb{F}_2^n .

³ $h(p) = -p \log_2(p) - (1-p) \log_2(1-p)$, for $p \in (0, 1)$, and $h(0) = h(1) = 0$.

⁴The fact that random linear codes are optimal for symmetric channels was shown in [Eli55].

⁵For the case of the BEC, [LMS⁺97] provides the first LDPC codes that are capacity-achieving, and further LDPC ensembles have been recently developed with spatial coupling [KRU11].

⁶The belief that RM codes achieve capacity is much older, but we did not trace back where it appears first.

constant, the most popular regime in coding theory.

The conjecture that RM codes achieve capacity has been experimentally “confirmed” in simulations [Ari08, MHU14]. Moreover, despite being extremely old, new interest in it resurged a few years ago with the advent of polar codes [Ari09]. To explain the connection between the two, as well as some of the technical problems arising in proving the results above, consider the following $2^m \times 2^m$ matrix E_m (for “evaluation”). Index the rows and columns by all possible m -bit vectors in \mathbb{F}_2^m in *lexicographic order*. Interpret the columns simply as points in \mathbb{F}_2^m , and the rows as monomials (where an m -bit string correspond to the monomial which is the product of variables in the positions containing a 1). Finally, $E_m(x, y)$ is the value of the monomial x on the point y (namely it is 1 if the set of 1’s in x is contained in the set of 1’s in y). Thus, every row of E_m is the truth table of one monomial. It is thus easy to see that the code $R(m, r)$ is simply the span of the top (or “high weight”) k rows of E_m , with $k = \binom{m}{\leq r}$; these are the truth tables of all degree $\leq r$ polynomials. In contrast, polar codes of the same rate are spanned by a different set of k rows, so they form a different subspace of polynomials. While the monomials indexing the polar code rows have no explicit description (so far), they can be computed efficiently for any k in $\text{poly}(n) = 2^{O(m)}$ time. It is somehow intuitively “better” to prefer higher weight rows to lower weight ones as the basis of the code (as the “chances of catching an error” seem higher). Given the amazing result that polar codes achieve capacity, this intuition seems to suggest that RM codes do so as well. In fact, experimental results in [MHU14] suggest that RM codes may outperform polar codes for the BEC and BSC with maximum-likelihood⁷ decoding.

Denoting by $E(m, r)$ the top submatrix of E_m with $k = \binom{m}{\leq r}$ rows, one can express some natural problems concerning it which are essential for our results. To obtain some of our results on achieving capacity for the erasure channel, we must understand the following two natural questions regarding $E(m, r)$. First, what is the largest number s so that s random columns of $E(m, r)$ are linearly independent with high probability. Second, what is the smallest number t such that t random columns have full row-rank. Capacity achieving for erasures means that $s = (1 - o(1))k$ and $t = (1 + o(1))k$, respectively. We prove that this is the case for small values of r . The second property gives directly the result for low-rate codes $\text{RM}(m, r)$, and the first implies the result for high-rate codes using a duality property of RM codes. Both results may be viewed from a learning theory perspective, showing that in these ranges of parameters any degree r polynomial in m variables can be uniquely interpolated with high probability from its values on the minimum possible number of random inputs.

For errors, further analysis is needed beyond the rank properties discussed above. From the parity-check matrix viewpoint, decoding errors is equivalent to solving (with high probability) an underdetermined system of equations. Recall that a linear code can be expressed as the null space of an $(n - k) \times n$ parity-check matrix H . If Z is a random error vector with about (or at most) s one’s corrupting a codeword, applying the parity-check matrix to the codeword yields $Y = HZ$, where the “syndrom” Y is of lower dimension $n - k$. Decoding random errors means reconstructing Z from Y with high probability, using the fact that Z is sparse (hence the connection with sparse recovery). Note however that this differs from compressed sensing, as Z is random and HZ is over $GF(2)$. It relates to randomness extraction in that a capacity achieving code should produce an output Y of dimension $m \approx nh(s/n)$ containing⁸ all the entropy of Z . Compared to the usual notion of randomness extraction, the challenge here is to extract with a very simple map H (seedless and linear), while the source Z is much more structured, i.e. it has i.i.d. components, compared to

⁷ML decoding looks for the most likely codeword. For the BEC, this requires inverting a matrix over $GF(2)$, whereas for the BSC, ML can be approximated by a successive list-decoding algorithm.

⁸See [Abb11] for further discussion on this.

sources in the traditional extractor settings.

Another extremely basic statistics of a code is its *weight distribution*, namely, approximately how many codewords have a given Hamming weight. Amazingly enough, very little was known about the weight distribution of Reed-Muller code until the recent breakthrough paper of [KLP12], who gave nearly tight bounds for constant degree polynomials for both. The results of [KLP12] also apply to list-decoding of RM codes, which was previously investigated in [GKZ08]. We need a sharpening of their upper bound for two of our results, which we prove by refining their method. The new bound is nearly tight not only for constant degree polynomials, but actually remains so even for degree r that is linear in m . We get a similar improvement for their bound on the list-size for list decoding of RM codes.

Summarizing, we study some very basic coding-theoretic questions regarding low-degree polynomials over $\text{GF}(2)$. We stress two central aspects which remain elusive. First, while proving the first results about parameters of RM codes which achieve capacity, the possibly most important range, when error rate is constant, seems completely beyond the reach of our techniques. Second, while our bounds for erasures immediately entails a (trivial) efficient algorithm to actually locate them, there are no known efficient algorithms for correcting random errors in the regimes we prove it is information theoretically possible. We hope that this paper will inspire further work on the subject, and we provide concrete open questions it suggests. We now turn to give more details on the problems, results and past related work.

1.2 Notation and terminology

Before presenting our results we need to introduce some notations and parameters. The following are used throughout the paper:

- For nonnegative integers $r \leq m$, $RM(m, r)$ denotes the Reed-Muller code whose codewords are the evaluation vectors of all multivariate polynomials of degree at most r on m Boolean variables. The maximal degree r is sometimes called the order of the code. The blocklength of the code is $n = 2^m$, the dimension $k = k(m, r) = \sum_{i=0}^r \binom{m}{i} \triangleq \binom{m}{\leq r}$, and the distance $d = d(m, r) = 2^{m-r}$. The code rate is given by $R = k(m, r)/n$.
- We use $E(m, r)$ to denote the “evaluation matrix” of parameters m, r , whose rows are indexed by all monomials of degree $\leq r$ on m Boolean variables, and whose columns are indexed by all vectors in \mathbb{F}_2^m . For $u \in \mathbb{F}_2^m$, we denote by u^r the column of $E(m, r)$ indexed by u , which is a k -dimensional vector, and for a subset of columns $U \subseteq \mathbb{F}_2^m$ we denote by U^r the corresponding submatrix of $E(m, r)$.
- A generator matrix for $RM(m, r)$ is given by $G(m, r) = E(m, r)$, and a parity-check matrix for $RM(m, r)$ is given by $H(m, r) = E(m, m - r - 1)$ (see Lemma 2.11).
- We associate with a subset $U \subseteq \mathbb{F}_2^m$ its characteristic vector $\mathbb{1}_U \in \{0, 1\}^n$. We often think of the vector $\mathbb{1}_U$ as denoting either an *erasure pattern* or an *error pattern*.

Finally, we use the following standard notations. $[n] = \{1, \dots, n\}$. The Hamming weight of $x \in \mathbb{F}_2^n$ is denoted $w(x) = |\{i \in [n] : x_i \neq 0\}|$ and the relative weight is $\text{wt}(x) = w(x)/n$. We use $B(n, s) = \{x \in \mathbb{F}_2^n : w(x) \leq s\}$ and $\partial B(n, s) = \{x \in \mathbb{F}_2^n : w(x) = \lceil s \rceil\}$. We use $\binom{[n]}{s}$ to denote the set of subsets of $[n]$ of cardinality s . Hence, for $S \in \binom{[n]}{s}$, $\mathbb{1}_S \in \partial B(n, s)$.

For a vector x of dimension n and subset S of n , we use $x[S]$ to denote the components of x indexed by S , and if X is matrix with n columns, we use $X[S]$ to denote the subset of columns indexed by S . In particular, $E(m, r)[U] = U^r$. When we need to be more explicit, for an $a \times b$

matrix A and $I \subseteq [a]$, we denote with $A_{I,\cdot}$, the matrix obtained by keeping only those rows indexed by I , and denote similarly $A_{\cdot,J}$ for $J \subseteq [b]$.

Channels, capacity and capacity-achieving codes

We next describe the channels that we will be working with, and provide formal definitions in Section 2. Throughout p will denote the corruption probability per coordinate. The Binary Erasure Channel (BEC) with parameter p acts on vectors $v \in \{0, 1\}^n$, by changing every coordinate to “?” with probability p . That is, after a message v is transmitted in the BEC the received message \hat{v} satisfies that for every coordinate i either $\hat{v}_i = v_i$ or $\hat{v}_i = \text{“?”}$ and $\Pr[\hat{v}_i = \text{“?”}] = p$. The Binary Symmetric Channel (BSC) with parameter p is flips the value of each coordinate with probability p . That is, after a message v is transmitted in the BSC the received message \hat{v} satisfies $\Pr[\hat{v}_i \neq v_i] = p$.

In fact, we will use a small variation on these channels; for corruption probability p we will fix the number of erasures/errors to $s = pn$. We note that by the Chernoff-Hoeffding bound (see e.g., [AS92]), the probability that more than $pn + \omega(\sqrt{pn})$ erasures/errors occur for independent Bernoulli choices is $o(1)$, and so we can restrict our attention to talking about a fixed number of erasures/errors. Thus, when we discuss s corruptions, we will take the corruption probability to be $p = s/n$. We refer to Section 2 for the details.

We now define the notions of “capacity-achieving” for the channels above. We consider $RM(m, r)$ where $r = r(m)$ typically depends on m . We say that $RM(m, r)$ can correct random erasures/errors, if it can correct the random erasures/errors with high probability when n tends to infinity. The goal is to recover from the largest amount of erasures/errors that is information-theoretically achievable. We note that while recovering from erasures, whenever possible, is always possible efficiently (by linear algebra), this need not be the case for recovery from errors. As we focus on the information theoretic limits, we allow maximum-likelihood (ML) decoding rule. Obtaining an efficient algorithm is a major open problem. Note that ML minimizes the error probability for equiprobable messages, hence if ML fails to decode the codewords with high probability, no other algorithms can succeed.

Recall that the capacity of a channel is the largest possible code rate at which we can recover (whp) from corruption probability p . This capacity is given by $1 - p$ for BEC erasures, and by $1 - h(p)$ for BSC errors. Namely, Shannon proved that for any code of rate R that allows to correct corruptions of probability p , then $R < 1 - p$ for the BEC and $R < 1 - h(p)$ for the BSC.

Achieving capacity means that R is close to the upper bound, say within $(1 + \varepsilon)$ factor of the optimal bounds above. For *fixed* corruption probabilities p and rates R in $(0, 1)$ this is easy to define (previous paragraph). However as we deal with very low or very high rates above, defining this needs a bit more care, and is described in the table below, and formally in Section 2. A code of rate R is ε -close to achieve capacity if it can correct from a corruption probability p that satisfies the bounds below⁹. It is capacity-achieving if it is ε -close to achieve capacity for all $\varepsilon > 0$.

	BEC	BSC
Low code-rate ($R \rightarrow 0$)	$p \geq 1 - R(1 + \varepsilon)$	$h(p) \geq 1 - R(1 + \varepsilon)$
High code-rate ($R \rightarrow 1$)	$p \geq (1 - R)(1 - \varepsilon)$	$h(p) \geq (1 - R)(1 - \varepsilon)$

⁹Note that for $R \rightarrow 0$, in the BEC we have $p \rightarrow 1$, while for the BSC we have $p \rightarrow \frac{1}{2}$. Also, we have stated the bounds thinking of R fixed and putting a requirement on p . One can equivalently fix p and require the code to correct a corruption probability p for a rate R that satisfies the bounds in the table.

1.3 Our results

We now state all our results, with approximate parameters, as the exact statements (given in the body of the paper) are somewhat technical. We divide this section to results on decoding from random erasures, then on weight distribution and list decoding, and finally decoding random errors. In brief, we investigate four cases: two regimes for the code rates (high and low rates) and two models (BEC and BSC). Besides for the BSC at high-rate, we obtain a capacity-achieving result for all other three cases. For the low-rate regimes, we obtain results for values of r up to the order of m .

1.3.1 Random erasures - the BEC channel

As mentioned earlier, some of the questions we study concerning properties of Reed-Muller codes can be captured by the following basic algebraic-geometric questions about evaluation vectors of low-degree monomials, namely, submatrices of $E(m, r)$. For any parameters $r \in [m]$ (the degree) and $s \in [n]$ (the size of the corrupted set U), we will study correcting random erasures and errors patterns of size s in $RM(m, r)$.

1. What is the largest s for which the submatrix U^r has full column-rank with high probability?
2. What is the smallest s for which the submatrix U^r has full row-rank with high probability?

More generally, we will be interested in characterizing sets U for which these properties hold. We note that for achieving capacity, s should be as close as possible to $\binom{m}{\leq r}$ (from below for the first question and from above for the second question). In other words, the matrix U^r should be as close to square as possible. Note that this would be achieved for the case where $E(m, r)$ is replaced by a random uniform matrix, so our goal in a sense is to show that $E(m, r)$ behaves like a random matrix with respect to these questions.

We obtain our decoding results for the BEC by providing answers to these questions for certain ranges of parameters. Our first theorem concerns Reed-Muller codes of low degree.

Theorem 1.1 (See Theorem 4.17). *Let $r = o(m)$. Then, If we pick uniformly at random a set U of $(1 + o(1)) \cdot \binom{m}{\leq r}$ columns of $E(m, r)$, then with probability $1 - o(1)$ the rows of this submatrix are linearly independent, i.e., U^r has full row-rank.*

As an immediate corollary we get that Reed-Muller codes of sub-linear degree achieve capacity for the BEC.

Theorem 1.2 (See Corollary 5.1). *For $r = o(m)$, $RM(m, r)$ achieves capacity for the BEC. More precisely, for every $\delta > 0$ and $\eta = O(1/\log(1/\delta))$ the following holds: For every $r \leq \eta m$, $RM(m, r)$ is δ -close to capacity for the BSC*

We obtain similar results in a broader range of parameters when the code is of high degree rather than low degree (i.e., the code has high rate rather than low rate).

Theorem 1.3 (See Theorem 4.5). *Let $r = O(\sqrt{m/\log m})$. If we pick uniformly at random a set U of $(1 - o(1)) \cdot \binom{m}{\leq r}$ columns of $E(m, r)$, then with probability $1 - o(1)$ they are linearly independent, i.e., the submatrix U^r has full column rank.*

Due to the duality between linear independent set of columns in $E(m, r)$ and spanning sets in the generating matrix of $RM(m, m - r - 1)$ (see Lemma 4.3) we get as corollary that Reed-Muller codes with the appropriate parameters achieve capacity for the BEC.

Theorem 1.4 (See Corollary 5.2). *For $m - r = O(\sqrt{m/\log m})$, $RM(m, r)$ is capacity-achieving on the BEC.*

1.3.2 Weight distribution and list decoding

Before moving to our results on random errors, we take a detour to discuss our results on the weight distribution of Reed-Muller codes as well as their list decoding properties. These are naturally important by themselves, and, furthermore, tight weight distribution bounds turns out to be crucial for achieving capacity for the BEC in Theorem 1.1 above, as well as for achieving capacity for the BSC in Theorem 1.7 below. Our bound extends an important recent result of Kaufman, Lovett and Porat on the weight-distribution of Reed-Muller codes [KLP12], using a simple variant of their technique. Kaufman et al. gave a bound that was tight for $r = O(1)$, but degrades as r grows. Our improvement extends this result to degrees $r = O(m)$. Denote with $W_{m,r}(\alpha)$ the number of codewords of $RM(m, r)$ that have at most α fraction of nonzero coordinates.

Theorem 1.5 (See Theorem 3.3). *Let $1 \leq \ell \leq r - 1 < m/4$ and $0 < \varepsilon \leq 1/2$. Then,*

$$W_{m,r}((1 - \varepsilon)2^{-\ell}) \leq (1/\varepsilon)^{O\left(\ell^4 \binom{m-\ell}{\leq r-\ell}\right)}.$$

As in the paper of [KLP12], almost the exact same proof as our proof of Theorem 1.5 yields a bound for list-decoding of Reed-Muller codes, for which we get similar improvements. Following [KLP12] we denote:

$$L_{m,r}(\alpha) = \max_{g: \mathbb{F}_2^m \rightarrow \mathbb{F}_2} |\{f \in RM(m, r) \mid \text{wt}(f - g) \leq \alpha\}|.$$

That is, $L_{m,r}(\alpha)$ denotes the maximal number of code words of $RM(m, r)$ in a hamming ball of radius $\alpha 2^m$. The bound concerns α of the form $(1 - \varepsilon)2^{-\ell}$ for $1 \leq \ell \leq r - 1$, and our main contribution is making the first factor in the exponent depend on ℓ (rather than on r in [KLP12]).

Theorem 1.6. *Let $1 \leq \ell \leq r - 1$ and $0 < \varepsilon \leq 1/2$. Then, if $r \leq m/4$ then*

$$L_{m,r}((1 - \varepsilon)2^{-\ell}) \leq (1/\varepsilon)^{O\left(\ell^4 \binom{m-\ell}{\leq r-\ell}\right)}.$$

1.3.3 Random errors - the BSC channel

We now return to discuss decoding from random errors. Our next result shows that Reed-Muller codes achieve capacity also for the case of random errors at the low rate regime. The proof of this result relies on Theorem 1.5.

Theorem 1.7 (See Theorem 6.1). *For $r = o(m)$, $RM(m, r)$ achieves capacity for the BSC. More precisely, for every $\delta > 0$ and $\eta = O(1/\log(1/\delta))$ the following holds: For every $r \leq \eta m$, $RM(m, r)$ is δ -close to capacity for the BSC.*

To obtain results about the behavior of high-rate Reed-Muller codes with respect to random errors we use a novel connection between robustness to errors and robustness to erasures in related Reed-Muller codes.

Theorem 1.8 (See Theorem 6.13). *If a set of columns U are linearly independent in $E(m, r)$ (namely, $RM(m, m - r - 1)$ can correct the erasure pattern $\mathbb{1}_U$), then the error pattern $\mathbb{1}_U$ can be corrected (i.e., it is uniquely decodable) in $RM(m, m - (2r + 2))$.*

Using Theorem 1.3 this gives a new result on correcting random errors in Reed-Muller codes.

Theorem 1.9 (See Theorem 6.2). *For $r = O(\sqrt{m/\log m})$, $RM(m, m - (2r + 2))$ can correct a random error pattern of weight $(1 - o(1)) \cdot \binom{m}{\leq r}$ with probability larger than $1 - o(1)$.*

While this result falls short of showing that Reed-Muller codes achieve capacity for the BSC in this parameter range, it does show that they can cope with many more errors than suggested by their minimum distance. Recall that the minimum distance of $R(m, m - (2r + 2))$ is 2^{2r+2} . Achieving capacity for this code means that it should be able to correct roughly $\binom{m}{2r}$ random errors. Instead we show that it can handle roughly $\binom{m}{\leq r}$ random errors, which is approximately the square root of the number of errors at capacity.

The proof of Theorem 1.8 reveals a more general phenomenon, that of reducing error correction to erasure correction. We prove that for any linear code C , of very high rate, there is another linear code C' of related high rate, so that if C can correct the *erasure* pattern $\mathbf{1}_U$ then C' can correct the *error* pattern $\mathbf{1}_U$. Furthermore C' is very simply defined from C . The decline in quality of C' relative to C is best explained in terms of the co-dimension (namely the number of linear constraints on the code, or equivalently the number of rows of its parity-check matrix). We prove that the co-dimension of C' is roughly the cube of the co-dimension of C . We now state this general theorem.

For a matrix H we denote by H^r the corresponding matrix that contains the evaluations of *all* columns of H by all degree $\leq r$ monomials (in an analogous way to the definition of U^r from U).

Theorem 1.10 (See Theorem 6.17). *If a set of columns U is linearly independent in a parity check matrix H , then the code that has H^3 as a parity check matrix can correct the error pattern $\mathbf{1}_U$.*

Note that applying this result as is to $E(m, r)$ would give a weaker statement than Theorem 1.8, in which $E(m, 2r + 1)$ would be replaced by $E(m, 3r)$. We conclude by showing that this result is tight, namely replacing 3 by 2 in the theorem above fails, even for RM codes.

Theorem 1.11 (See Section 6.2.5). *There are subsets of columns U that are linearly independent in $E(m, 1)$, but such that the patterns $\mathbf{1}_U$ are not uniquely decodable in $E(m, 2)$.*

1.4 Proof techniques

Although the statements of Theorems 1.3 and 1.1 sound very similar, their proofs are very different. We first explain the ideas behind the proofs of these two theorems and then give details for the proofs of Theorems 1.5, 1.7, 1.8 and 1.10.

Proof of Theorem 1.3 The proof of Theorem 1.3 relies on estimating the size of varieties (sets of common zeros) of linear subspaces of degree r polynomials. Here is a high level sketch.

Recall that we have to show that if we pick a random set of points $U \subset \mathbb{F}_2^m$, of size $(1 - o(1)) \cdot \binom{m}{\leq r}$, and with each point associate its degree- r evaluation vector, then with high probability these vectors are linearly independent. While proving this is simple when considered over large fields, it is quite non-trivial over very small fields. We are able to prove that this property holds for degrees r up to (roughly) $\sqrt{m/\log m}$. It is a very interesting question to extend this to larger degrees as well.

To prove that a random set U of appropriate size gives rise to linearly independent evaluation vectors we consider the question of what it takes for a new point to generate an evaluation vector that is linearly independent of all previously generated vectors. As we prove, this boils down to understanding what is the probability that a random point is a common zero of all degree r polynomials, in a certain linear space of polynomials defined by the previously picked points. If this set of common zeros is small, then the success probability (i.e., the probability that a new point will yield an independent evaluation vector) is high, and we can iterate this argument.

To bound the number of common zeros we yet again move to a dual question. Notice that if a set of K linearly independent polynomials of degree r vanishes on a set of points V , then there are at

most $\binom{m}{\leq r} - |K|$ linearly independent degree r polynomials that are defined over V . In view of this, the way to prove that a given set of polynomials does not have too many common zeros is to show that any large set of points (in our case, the set of common zeros) has many linearly independent degree r polynomials that are defined over it. We give two different proofs of this fact. The first uses a hashing argument; if V is large then after some linear transformation it supports many different degree r *monomials*. The second relies on a somewhat tighter bound that was obtained by Wei [Wei91], who studied the generalized Hamming weight of Reed-Muller codes. While Wei’s result gives slightly tighter bounds compared to the hashing argument, we find the latter argument more transparent.

Proof of Theorem 1.1 To prove Theorem 1.1 we first observe that a set of columns U (in $E(m, r)$) spans the entire row-space if and only if there is no linear combination of the rows of $E(m, r)$ that is supported on the complementary set $U^c = \mathbb{F}_2^m \setminus U$. As linear combinations of rows correspond to “truth-tables” of degree r polynomials, this boils down to proving that, with high probability, no nonzero degree r polynomial vanishes on all points in U . For each such polynomial, if we know its weight (the number of nonzero values it takes), this is a simple calculation, and the hope is to use a union bound over all polynomials. To this end, we can partition the codewords to dyadic intervals according to their weights, carry out this calculation and union bound the codewords in each interval and then combine the results. For this plan to work we need a good estimate of the number of codewords in each dyadic interval, which is what Theorem 1.5 gives.

Proof of Theorem 1.5 As mentioned earlier, this theorem improves upon a beautiful result of Kaufman, Lovett and Porat [KLP12]. Our proof is closely related to their proof. Roughly, what they show is that any small weight codeword, i.e., a degree r polynomial with very few non-zero values, can be well approximated by a “few” partial derivatives. Namely, there is a function that when applied to a few lower degree polynomials, agrees with the original polynomial on most of the inputs. Here “few” depends on the degree r , the weight and (crucially for us) the quality of the approximation. Kaufman et al. then pick an approximation quality parameter that guarantees that the approximating function can be close to at most *one* polynomial of degree r . Then, counting the number of possible approximating functions they obtain their bound. The cost is that such a small approximation parameter blows the number of “few” derivatives that are required. We diverge from their proof in that we choose a much larger approximation quality parameter, but rather allow each approximating function to be close to *many* degree r polynomials. The point is that, by the triangle inequality, all these polynomials are close to each other, and so subtracting one of them from any other still yield polynomials of very small weight. Thus, we can use induction on weight to bound their number, obtaining a better bound on the number of polynomials of a given weight.

Proof of Theorem 1.7 Here we use a coarse upper-bound on the error probability, as for the proof that a random linear code achieves capacity, and show that the argument still holds for RM codes. To prove that a random code can, w.h.p., uniquely decode an error pattern $\mathbb{1}_U$ of weight w we basically wish to show that for no other error pattern $\mathbb{1}_V$, of weight w , the vector $\mathbb{1}_U \oplus \mathbb{1}_V$ is a code word (as then both error patterns will have the same syndrome). Stated differently, we want to count how many different ways are there to represent a codeword as a sum of two vectors of weight at most w . This counting depends very much on the weight of the codeword that we wish to split. In random linear codes weights are very concentrated around $n/2$, which makes a union bound easy. Reed-Muller codes however have many more codewords of smaller weights, and the argument depends precisely on how many. Once again Theorem 1.5 comes to the rescue and

enables us to make this delicate calculation for each (relevant) dyadic interval of weights. Here too our improvement of [KLP12] is essential.

Proofs of Theorems 1.8, 1.9 and 1.10 Consider an erasure pattern $\mathbb{1}_U$ such that the corresponding set of degree- r evaluation vectors, U^r , is linearly independent. Namely the columns indexed by U in $E(m, r)$ are linearly independent. We would like to prove that $\mathbb{1}_U$ is uniquely decodable from its syndrome under $H(m, m - 2r - 2) = E(m, 2r + 1)$. We actually prove that if $\mathbb{1}_V$ is another erasure pattern, which has the same syndrome under $H(m, m - 2r - 2)$, then $U = V$. The proof may be viewed as a reconstruction (albeit inefficient) of the set U from its syndrome. Here is a high level description of our argument that different (linearly independent) sets of erasure patterns give rise to different syndromes.

We first prove this property for the case $r = 1$ (details below). This immediately implies Theorem 1.10 as every parity check matrix of any linear code is a submatrix of $E(m, 1)$ for some m . This is a general reduction from the problem of recovering from errors to that of recovering from erasures (in a related code). As a special case, it also implies that for any r , $H(m, m - 3r - 1)$ uniquely decodes any error pattern $\mathbb{1}_U$ such that the columns indexed by elements of U in $E(m, r) = H(m, m - r - 1)$ are linearly independent. We then slightly refine the argument for larger degree r to replace $H(m, m - 3r - 1)$ above by $H(m, m - 2r - 2)$, which gives Theorem 1.8.

For the case $r = 1$, the proof divides to two logical steps. In the first part we prove that the columns of V must span the same space as the columns of U . This requires only the submatrix $E(m, 2)$, namely at pairs of coordinates in each point (degree 2 monomials). In the second part we use this property to actually identify each vector of U inside V . This already requires looking at the full matrix $E(m, 3)$, namely at triples of coordinates.

It is interesting that going to triples of coordinate is essential for $r = 1$ (and so this result are tight). We prove that even if the columns of U are linearly independent, then there can be a different set V that has the same syndrome in $E(m, 2)$. This result is given in Section 6.2.5. We do not know what is the right bound for general r .

1.5 Related literature

Recovery from random corruptions

Besides the conjectures mentioned in the introduction that RM codes achieve capacity, results fall short of that for all but very special cases. We are not familiar of works correcting random erasures. Several papers have considered the quality of RM codes for correcting random errors when using specific algorithms, focusing mainly on efficient algorithms. In [Kri70], the majority logic algorithm [Ree54] is shown to succeed in recovering all but a vanishing fraction of error patterns of weight up to $d \log(d)/4$, where $d = 2^{m-r}$ is the code distance, requiring however a positive rate $R > 0$. This was later improved to weights up to $d \log(d)/2$ in [DS06]. We note that for a fixed rate $0 < R < 1$, this is roughly \sqrt{n} , whereas to achieve capacity one should correct $\Omega(n)$ erasures/errors.

A line of work by Dumer [Dum04, DS06, Dum06] based on recursive algorithms (that exploits the recursive structure of RM codes), obtains results mainly for low-rate regimes. In [Dum04], it is shown that for a fixed order r , i.e., for $k(m, r) = \Theta(m^r)$, an algorithm of complexity $O(n \log(n))$ can correct most error patterns of weight up to $n(1/2 - \varepsilon)$ given that ε exceeds $n^{-1/2^r}$. In [DS06], this is improved to errors of weight up to $n/2(1 - (4m/d)^{1/2^r})$, requiring that $r/\log(m) \rightarrow 0$. Further, [Dum06] shows that most error patterns of weight up to $n/2(1 - (4m/d)^{1/2^r})$ can be recovered in the regime where $\log(m)/(m - r) \rightarrow 0$.

Note that while the previous results rely on efficient decoding algorithms, they are far from being capacity-achieving. Concerning maximum-likelihood decoding, known algorithms are of

exponential complexity in the blocklength, besides for special cases such as $r = 1$ or $r = m - 2$ [AL04]. In [VMS92], it is shown that RM codes of fixed order r can decode most error patterns of weight up to $n/2(1 - \sqrt{c(2^r - 1)m^r/nr!})$, where $c > \ln(4)$. However, this does not provide a capacity-achieving result, which would require decoding most error patterns of weight approaching $n/2(1 - \sqrt{\ln(4)m^r/nr!})$, i.e., [VMS92] has an extra $\sqrt{2^r - 1}$ factor.

For the special case of $r = 1, 2$ (i.e., the generator matrix has only vectors of weights n , $n/2$ and $n/4$), [HKL05] shows that RM codes are capacity-achieving. For $r \geq 3$, the problem is left open.

Weight enumeration

The weight enumerator (how many codewords are of any given weight) of $RM(m, 2)$ was characterized in [SB70]. For $RM(m, 3)$, a complete characterization of the weight enumerator is still missing. The number of codewords of minimal weight is known for any r , and corresponds to the number of $(m - r)$ -flats in the affine geometry $AG(m, 2)$ [MS77]. In [KT70], the weight enumerator of RM codes is characterized for codewords of weight up to twice the minimum distance, later improved to 2.5 the minimum distance in [KTA76].

For long, [KTA76] remained the largest range for which the weight enumerator was characterized, until [KLP12] managed to breakthrough the 2.5 barrier and obtained bounds for all distances in the regime of small r . The results of [KLP12] is given in Theorem 3.1.

1.6 Organization

The paper is organized as follows. We first discuss the model of random erasures and errors (Section 2) and then give a quick introduction to Reed-Muller codes (Section 2.3). In section 3 we prove Theorem 1.5 on the weight distribution of RM codes. In Section 4 we give answers to the two questions on sub matrices of $E(m, r)$, when r is small. In Section 5 we use the result obtained thus far to obtain our results for the BEC (Theorems 1.1 and 1.3). In Section 6 we give our results for the BSC. Finally, in Section 7, we discuss some intriguing future directions and open problems which our work raises.

2 Preliminaries

In this section, we review basic concepts about linear codes and their capability of correcting random corruptions, as well as Reed-Muller codes.

2.1 Basic coding definitions

Recall that for a binary linear code $C \subseteq \mathbb{F}_2^n$ of blocklength n , if k denotes the dimension of a code, i.e., $k = \log_2 |C|$, a (non-redundant) generator matrix G has dimension $k \times n$, a (non-redundant) parity-check matrix H has dimension $(n - k) \times n$, and $C = \text{Im}(G) = \ker(H)$.

In the worst-case model, the distance of the code determines exactly how many erasures and errors can be corrected, with the following equivalent statements for the generator and parity-check matrices:

- C has distance d ,
- C allows to correct $d - 1$ erasures,
- C allows to correct $\lfloor (d - 1)/2 \rfloor$ errors,
- any $d - 1$ columns of H are linearly independent,

- any $n - d + 1$ rows of G have full span.

Two fundamental problems in worst-case coding theory is to determine the largest dimension of a code that has distance at least d , for a fixed d , and to construct explicit codes achieving the optimal dimension. None of these questions are solved in general, nor in the asymptotic regime of n tending to infinity with $d = \alpha n$, and $\alpha \in (0, 1/2)$. A random linear code achieves a dimension of $n(1 - h(\alpha)) + o(n)$, the Gilbert-Varshamov bound, but this bound has not been proved to be tight nor has it been improved asymptotically since 1957. Further, no explicit construction is known to achieve this bound.

In this paper, we are interested in random erasures and errors, and in correcting them “with high probability.” This changes the requirements on the generator and parity-check matrices. For erasures, it simple requires linear independence of the columns “with high probability” (see Section 2.2), while the requirement is more subtle for errors. Yet, in the probabilistic setting, random codes can be proved to achieve the optimal tradeoffs (e.g., code rate vs. erasure rate, or code rate vs. error rate) that are known for both erasures and errors (as special cases of Shannon’s theorem). Explicit constructions of codes that achieve the optimal tradeoffs are also known, e.g., polar codes [Ari09], and this paper investigates RM codes as such candidates. We now provide the formal models and results.

We mainly work in this paper with “uniform” models for erasures and errors, but we sometimes interpret the results for “i.i.d.” models (namely, the BEC and BSC channels). To formally connect these, we define first a unified probabilistic model. We start with erasures. We restrict ourselves to linear codes, although the definitions extend naturally to non-linear codes.

Definition 2.1. *A sequence of linear codes $\{C_n\}_{n \geq 1}$ of blocklength n allows to correct random erasures from a sequence of erasure distributions $\{\mu_n\}_{n \geq 1}$ on \mathbb{F}_2^n , if for $Z \sim \mu_n$ and $S_Z = \{i \in [n] : Z_i = 1\}$ (the erasure pattern),*

$$\Pr\{\exists x, y \in C_n : x \neq y, x[S_Z^c] = y[S_Z^c]\} \rightarrow 0, \quad \text{as } n \rightarrow \infty,$$

i.e., the probability of drawing erasures at locations that can confuse different codewords is vanishing.

Notice that $x[S_Z^c] = y[S_Z^c]$ if and only if we cannot correct erasures on coordinates S_Z for neither x nor y . We now present a unified model for errors.

Definition 2.2. *A sequence of linear codes of length n and parity-check matrix $\{H_n\}_{n \geq 1}$ allows to correct random errors from a sequence of error distributions $\{\mu_n\}_{n \geq 1}$ on \mathbb{F}_2^n if for $Z \sim \mu_n$,*

$$\Pr\{\exists z' \in \mathbb{F}_2^n : z' \neq Z, H_n z' = H_n Z, \mu_n(z') \geq \mu_n(Z)\} \rightarrow 0, \quad \text{as } n \rightarrow \infty, \quad (1)$$

i.e., the probability of drawing an error pattern Z for which there exists another error pattern z' that has the same syndrome as Z and is more likely than Z is vanishing.

Note that (1) is the requirement that the error probability of the maximum likelihood (ML) decoder vanishes. Since ML minimizes the error probability for equiprobable codewords, (1) is necessary for any decoder to succeed with high probability.

Remark 1. We next drop the term “sequence of” and the subscripts n , and simply say that a code C of blocklength n allows to correct random erasures/errors in specified models. The parameters introduced may also depend on n without being explicitly mentioned.

We now introduce the uniform and i.i.d. models.

Definition 2.3.

- (i) A linear code of blocklength n allows to correct $s = s_n$ random erasures (resp. errors) if it can correct them from the uniform erasure (resp. error) distribution U_s , i.e., the uniform probability distribution on $\partial B(n, s) = \{z \in \mathbb{F}_2^n : w(z) = \lceil s \rceil\}$.
- (ii) A linear code of blocklength n allows to correct erasures (resp. errors) for the $BEC(p)$ channel (resp. $BSC(p)$ channel), where $p = p_n$, if it can correct the distribution B_p , where B_p is the i.i.d. distribution¹⁰ on \mathbb{F}_2^n with Bernoulli(p) marginal.

Note that for $\mu_n = U_s$, i.e., the uniform distribution over $\partial B(n, s)$, the above definition reduces¹¹ to

$$\frac{|\{z \in \partial B(n, s) : \exists z' \text{ s.t. } z \neq z', Hz = Hz'\}|}{\binom{n}{s}} \rightarrow 0, \quad \text{as } n \rightarrow \infty,$$

i.e., the fraction of bad error patterns, which have non-unique syndrome, is vanishing.

The following Lemma¹² follows from standard probabilistic arguments.

Lemma 2.4.

- (i) If a linear code can correct $s = s_n$ random erasures (resp. errors), then it can correct erasures (resp. errors) from the $BEC((s - \omega(\sqrt{s}))/n)$ channel (resp. $BSC((s - \omega(\sqrt{s}))/n)$ channel).
- (ii) If a linear code can correct erasures (resp. errors) from the $BEC(p)$ channel (resp. $BSC(p)$ channel), then it can correct $np - \omega(\sqrt{np})$ random erasures (resp. errors).

We now define the notions of capacity-achieving. Since in the rest of the paper typically considers codes at a given rate, and investigate how many corruptions they can correct, the definitions are stated accordingly. Note that what follows is simply a restatement of Shannon's theorems for erasures and errors, namely that a code C of rate $R = (\log_2(|C|))/n$ correcting a corruption probability p must satisfy $R < 1 - p$ for erasures and $R < 1 - H(p)$ for errors. However, since we consider code rates that tend to 0 and 1, the requirements are broken down in various cases to prevent meaningless statements.

Definition 2.5. A code is capacity-achieving (or achieves capacity) if it is ε -close to capacity for all $\varepsilon > 0$. We now define the notion of ε -close to capacity in the four configurations:

- A linear code C of rate $R = o(1)$ is ε -close to capacity-achieving for erasures or for the BEC if it can correct np random erasures for a p satisfying

$$p \geq 1 - R(1 + \varepsilon).$$

- A linear code C of rate $R = o(1)$ is ε -close to achieving capacity for errors or for the BSC if it can correct np random errors for a p that satisfies

$$h(p) \geq 1 - R(1 + \varepsilon),$$

where $h(p) = -p \log_2(p) - (1 - p) \log_2(1 - p)$ is the entropy function.

¹⁰This means the product distribution with identical marginals.

¹¹We define $\binom{n}{s}$ as $\binom{n}{\lceil s \rceil}$ for a non-integer s .

¹²The statements are relevant for s_n or np_n that are $\omega(1)$.

- A linear code C of rate $R = 1 - o(1)$ is ε -close to achieving capacity for erasures or for the BEC if it can correct np random erasures for a p that satisfies

$$p \geq (1 - R)(1 - \varepsilon).$$

- A linear code C of rate $R = 1 - o(1)$ is ε -close to achieving capacity for erasures or for the BSC if it can correct np random erasures for a p that satisfies

$$h(p) \geq (1 - R)(1 - \varepsilon).$$

Note that previous definition leads to the same notion of capacity for the uniform and i.i.d. models in view of Lemma 2.4.

2.2 Equivalent requirements for probabilistic erasures

In this section, we show the following basic results: a code can correct s random erasures with high probability (whp), if a random subset of s columns in its parity-check matrix are linearly independent whp, or if a random subset of $n - s$ rows in its generator matrix have full-span whp.

First note the following algebraic equivalence, which simply states that a bad erasure pattern, one that can confuse different codewords, corresponds to a subset of rows of the generator matrix that is not full-span (i.e., not invertible).

Lemma 2.6. For an $n \times k$ matrix¹³ G and for $S \subseteq [n]$, let $G_{S,\cdot}$ denote the subset of rows of G indexed by S . Then the set of bad erasure patterns is given by

$$\left\{ S \in \binom{[n]}{s} : \exists x, y \in \ker(H), x \neq y, x[S^c] = y[S^c] \right\} \equiv \{ D \in \partial B(n, s) : \text{rank}(G_{D^c,\cdot}) < k \},$$

where $\text{rank}(G_{D^c,\cdot}) < k$ means that the columns of $G_{D^c,\cdot}$ are linearly dependent (i.e., the rows have full span).

Proof. We have

$$\begin{aligned} & \left\{ S \in \binom{[n]}{s} : \exists x, y \in \text{Im}(G), x \neq y, x[S^c] = y[S^c] \right\} \\ & \equiv \left\{ S \in \binom{[n]}{s} : \exists v \in \text{Im}(G) \text{ s.t. } v[S^c] = 0, v \neq 0 \right\} \\ & \equiv \left\{ S \in \binom{[n]}{s} : \text{rank}(G_{S^c,\cdot}) < k \right\}. \end{aligned}$$

□

Corollary 2.7. For an $n \times k$ matrix G and $s \in [n]$, denote by $G_{s,\cdot}$ the random sub-matrix of G obtained by selecting s rows uniformly at random. Then, the code with generator matrix G can correct s erasures if and only if

$$\Pr\{\text{rank}(G_{n-s,\cdot}) = k\} \rightarrow 1, \quad \text{as } n \rightarrow \infty.$$

¹³We assume that G has full column-rank, i.e., the generator matrix is non-redundant.

We now switch to the parity-check matrix interpretation. The following lemma shows that a bad erasure pattern for a code $C = \ker(H)$ can be identified as a subset of linear dependent columns in the parity-check matrix.

Lemma 2.8. *For a matrix H with n columns, for $S \subseteq [n]$, and $H[S]$ the subset of columns of H indexed by S , then the set of bad erasure patterns is given by*

$$\left\{ S \in \binom{[n]}{s} : \exists x, y \in \ker(H), x \neq y, x[S^c] = y[S^c] \right\} \equiv \left\{ S \in \binom{[n]}{s} : \text{rk}(H[S]) < s \right\},$$

where $\text{rk}(H[S]) < s$ simply means the columns of $H[S]$ are linearly dependent.

Proof of Lemma 2.8. Let

$$\begin{aligned} \text{BadSet} &:= \left\{ D \in \binom{[n]}{s} : \text{rk}(H[D]) < s \right\}, \\ \text{BadEra} &:= \left\{ S \in \binom{[n]}{s} : \exists x, y \in \ker(H), x \neq y, x[S^c] = y[S^c] \right\}, \end{aligned}$$

denote respectively the set of bad sets for which the columns of H do not have full rank and the set of bad erasure patterns that can confuse codewords in $\ker H$. Since the code is linear,

$$\text{BadEra} = \left\{ S \in \binom{[n]}{s} : \exists v \in \ker(H), v \neq 0, v[S^c] = 0 \right\}.$$

Hence for any $S \in \text{BadEra}$, there exists $v \in \ker(H)$, such that $v \neq 0$ and $\text{supp}(v) \subseteq S$, and the columns of H indexed by S are not full rank. Conversely, if $D \in \text{BadSet}$, D contains a subset V such that $V \neq \emptyset$ and $\sum_{i \in V} H[j] = 0$, hence $D \in \text{BadEra}$ (taking v as the indicator vector of V). \square

Corollary 2.9. *For a matrix H with n columns and $s \in [n]$, denote by $H[s]$ the random sub-matrix of H obtained by selecting s columns uniformly at random. Then, the code $\ker(H)$ can correct s random erasures if and only if*

$$\Pr\{\text{rk}(H[s]) = s\} \rightarrow 1, \quad \text{as } n \rightarrow \infty.$$

In other words, correcting s random erasures is equivalent to asking that a random subset of s columns in the parity-check matrix H is full rank whp.

While the requirement to correct probabilistic erasures (Corollaries 2.9 and 2.7) is similar to the requirement for the worst-case model but “with high probability,” the situation is more subtle in the case of errors. Note that for a code C with parity-check matrix H , the set of bad error patterns are the ones which lead to a non-unique syndrome, i.e.

$$\begin{aligned} &\{z \in \partial B(n, s) : \exists z' \in \partial B(n, s) \text{ s.t. } z \neq z', Hz = Hz'\} \\ &\equiv \{z \in \partial B(n, s) : \exists z' \in \partial B(n, s) \text{ s.t. } z \neq z', z + z' \in C\}. \end{aligned}$$

In other words, the set of bad error patterns are obtained by taking the set of codewords and splitting the codewords into elements of weight s . It is of course enough to consider the codewords of weight at most $2s$. However, even if the probability of drawing a codeword of weight at most $2s$ is vanishing, it does not follow that the probability of having a bad vector of weight s is also vanishing. There are multiple ways to split a codeword in vectors of weight s , and these lead to overlapping sets of vectors. Hence, the probability of a bad error pattern depends on the structure of H beyond the probability of having dependent columns.

2.3 Basic properties of Reed-Muller codes

The goal of this section is to revise the duality property of RM codes, which we use in this paper. One of the simplest way to understand this property is via the recursive structure of RM codes, mentioned below. We start by repeating the formal definition of RM codes via polynomials.

Definition 2.10. *Let m, r be two positive integers with $r \leq m$, and let $n = 2^m$. The Reed-Muller code of parameters m and r is defined by the set of codewords*

$$RM(m, r) = \{(f(a_0), \dots, f(a_{n-1})) : f \in \mathbb{P}(m, r)\},$$

where $\mathbb{P}(m, r)$ is the set of m -variate polynomials of degree at most r on \mathbb{F}_2 , and a_0, \dots, a_{n-1} are all the elements of \mathbb{F}_2^m .

In particular, the matrix $E(m, r)$ that contains only the evaluations of the monomials of degree at most r clearly defines a generator matrix for $RM(m, r)$. Formally, one should take the transpose $E(m, r)^t$ to obtain a generator matrix of $RM(m, r)$ that has dimension $n \times k$ (and not $k \times n$), where k is the dimension of the code (as usually assumed in coding theory and as in previous section). The duality property says that $E(\cdot, \cdot)$ can also be used to obtain a parity-check matrix of RM codes, as follows.

Lemma 2.11. *[Duality of RM codes] $E(m, m - r - 1)$ is a parity-check matrix for $RM(m, r)$, or equivalently, $E(m, r)$ is a parity-check matrix for $RM(m, m - r - 1)$.*

To show this result, note that RM codes can be defined recursively as follows. Instead of displaying the rows of the generator matrix by increasing order of the monomial degrees, consider the lexicographic order. For example, $RM(3, 3)$ is generated by :

$$\begin{array}{l} \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \\ 1 \quad \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{pmatrix} \\ x_1 \quad \begin{pmatrix} 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{pmatrix} \\ x_2 \quad \begin{pmatrix} 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \end{pmatrix} \\ x_1 x_2 \quad \begin{pmatrix} 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \end{pmatrix} \\ x_3 \quad \begin{pmatrix} 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{pmatrix} \\ x_1 x_3 \quad \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 \end{pmatrix} \\ x_2 x_3 \quad \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \end{pmatrix} \\ x_1 x_2 x_3 \quad \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \end{array}$$

Notice the order of x_3 and $x_1 x_2$ in the above. With that order, the matrix is the tensor product of $\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$ with itself 3 times. In fact, RM codes can equivalently be defined in terms of tensor products.

Definition 2.12. *For an integer $m \geq 0$, define*

$$G(m) = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}^{\otimes m},$$

with $G(0) = 1$. For $0 \leq r \leq m$, define $G(m, r)$ as the sub-matrix of $G(m)$ obtained by keeping the rows with weight more or equal to 2^{m-r} .

Note that $G(m, r)$ is simply a permutation of the rows of $E(m, r)$, hence it is also a generator matrix for $RM(m, r)$. Moreover, it can be constructed recursively as follows:

$$G(m, r) = \begin{pmatrix} G(m-1, r) & G(m-1, r) \\ 0 & G(m-1, r-1) \end{pmatrix}.$$

The polynomial interpretation of this recursion is simply the fact that a m -variate polynomial f of degree at most r can be expressed as

$$f(x_1, \dots, x_m) = f_1(x_1, \dots, x_{m-1}) + x_m f_2(x_1, \dots, x_{m-1}),$$

where f_1 and f_2 are m -variate polynomials of degrees at most r and $r-1$ respectively. With this recursion, the duality property (Lemma 2.11) (as well as the fact that the distance of $RM(m, r)$ is 2^{m-r}) are directly proved by induction. We refer to [MS77] for complete proofs.

3 Weight distribution of Reed-Muller codes

In this section we study the weight distribution of Reed-Muller codes. Our analysis is based on the technique of Kaufman, Lovett and Porat [KLP12]. We start with some high level intuition. Naturally, one expects that most codewords of $RM(m, r)$ (or any linear code, for that matter) to have weight around $n/2 = (2^m)/2$. A trivial upper bound on the number of codewords having such weight (or larger) is the total number of codewords, i.e., $2^{\binom{m}{r}}$. The question is thus how does this number changes when we consider smaller weights. Specifically, what is the number of codewords that have weight at most $2^{m-\ell}$ for some parameter ℓ . If we denote this number with $2^{c(m, r, \ell) \cdot \binom{m}{r}}$, then we are asking for the value of the term $c(m, r, \ell)$. A trivial lower bound on the number of such codewords is $2^{m\ell + \binom{m-\ell}{r-\ell}}$, which is obtained by counting all polynomials of degree- r that are divisible by ℓ linear functions. If this was tight, then $c(m, r, \ell) \approx 1$, which suggests that the number of such polynomials grows roughly like the number of degree $r-\ell$ polynomials on $m-\ell$ variables. Kaufman et al. proved that indeed this number is essentially the right answer for constant r . More precisely, they proved that $c(m, r, \ell) = O(r^2)$. Our contribution is replacing this estimate with roughly $c(m, r, \ell) = O(\ell^4)$. This change is most significant when ℓ is very small compared to r , e.g., when considering the number of words of weight e.g. roughly $n/4$ (so ℓ is 2) and when r is large, e.g. $r = \Omega(m)$. This improvement turns out to be critical for two of our results on achieving capacity – for erasure in low rates and errors in high rate. It remains open if one can prove that $c(m, r, \ell) = O(1)$, namely is a constant independent of all parameters.

We start by giving the high level view of the proof of [KLP12] and then explain how to improve their analysis. We first introduce some notation.

For a function $f : \mathbb{F}_2^m \rightarrow \mathbb{F}_2$ (equivalently, a word $f \in \mathbb{F}_2^n$) we denote by $\text{wt}(f)$ the relative (Hamming) weight of f , i.e.,

$$\text{wt}(f) = \frac{1}{2^m} |\{v \in \mathbb{F}_2^m \mid f(v) \neq 0\}|.$$

The cumulative weight distribution of $RM(m, r)$ at a relative weight $0 \leq \alpha \leq 1$, denoted $W_{m, r}(\alpha)$, is the number of codewords of $RM(m, r)$ whose relative weight is at most α ,

$$W_{m, r}(\alpha) \triangleq |\{f \in RM(m, r) \mid \text{wt}(f) \leq \alpha\}|.$$

The main theorem of [KLP12] roughly states that the number of code words of $RM(m, r)$ of relative weight at most $2^{-\ell}$ is roughly¹⁴ $\exp\left(r^2 \binom{m}{\leq r} \cdot \left(\frac{r}{m-r}\right)^\ell\right)$. I.e., the number of codewords of relative weight smaller than $1/2$ is significantly smaller than the number of words of relative weight $1/2$. We next give a slightly informal statement of the main theorem of [KLP12].¹⁵

Theorem 3.1 (Theorem 3.1 of [KLP12]). *Let $1 \leq \ell \leq r - 1$ and $0 < \varepsilon \leq 1/2$ be such that $2^{-r} \leq (1 - \varepsilon)2^{-\ell} < \frac{1}{2}$. Then*

$$W_{m,r}((1 - \varepsilon)2^{-\ell}) \leq (1/\varepsilon)^{O(r^2 \binom{m}{\leq r-\ell})}.$$

We next explain the main lemma used to prove Theorem 3.1. First we introduce the notion of a discrete partial derivative.

The discrete derivative of $f : \mathbb{F}_2^m \rightarrow \mathbb{F}_2$ in direction $y \in \mathbb{F}_2^m$ at point x is

$$\Delta_y f(x) \triangleq f(x + y) + f(x).$$

It is clear that $\Delta_y f(x) = \Delta_y f(x + y)$, so in particular, $\Delta_y \Delta_y f(x) = 0$. Thus, the function $\Delta_y f(\cdot)$ is determined by its value on the quotient space $\mathbb{F}_2^m / \langle y \rangle$, where for a set of vectors V , $\langle V \rangle$ denotes the space spanned by the vectors in V . It is a straight forward observation that if f is a polynomial of degree at most r then $\Delta_y f(\cdot)$ is a polynomial of degree at most $r - 1$. Similarly, the ℓ 'th iterated derivative of f in direction $Y = (y_1, \dots, y_\ell)$ and point x is

$$\Delta_Y f(x) \triangleq \Delta_{y_1} \Delta_{y_2} \dots \Delta_{y_\ell} f(x) = \sum_{I \subseteq [\ell]} f(x + \sum_{i \in I} y_i).$$

It is easy to show that $\Delta_Y f(x)$ does not depend on the order in which we take the derivatives.

We are now ready to state the main lemma of [KLP12].

Lemma 3.2 (Lemma 2.1 of [KLP12]). *Let $f : \mathbb{F}_2^m$ be such that $wt(f) \leq (1 - \varepsilon)2^{-\ell}$, for $0 < \varepsilon < 1$. Let $\delta > 0$ be an approximation parameter. There exists a universal algorithm \mathcal{A} (which does not depend on f) with the following properties:*

1. \mathcal{A} has two inputs: $x \in \mathbb{F}_2^m$ and $(Y_1, \dots, Y_t) \in (\mathbb{F}_2^m)^\ell$.
2. \mathcal{A} has oracle access to the ℓ 'th derivatives $\Delta_{Y_1} f(\cdot), \dots, \Delta_{Y_t} f(\cdot)$.

Then, for $t = c(\log(1/\delta) \log(1/\varepsilon) + \log^2(1/\delta))$, for some absolute constant c that does not depend on any of the parameters, there exists a setting for Y_1, \dots, Y_t such that

$$\Pr_{x \in \mathbb{F}_2^m} [\mathcal{A}(x, (Y_1, \dots, Y_t); \Delta_{Y_1} f(\cdot), \dots, \Delta_{Y_t} f(\cdot)) = f(x)] \geq 1 - \delta,$$

where $\Pr_{x \in \mathbb{F}_2^m}$ means that x is uniformly drawn in \mathbb{F}_2^m .

In other words, what the lemma shows is that if f has relatively low weight, then given an appropriate set of $O(\log(1/\delta) \log(1/\varepsilon) + \log^2(1/\delta))$ many ℓ -th derivatives of f , one can determine the value of f on most inputs. When f is a degree- r polynomial, its derivatives are degree $r - \ell$ polynomials and thus the lemma lets us approximate f well by lower degree polynomials.

¹⁴We use $\exp(x)$ instead of e^x .

¹⁵Kaufmann et al. also give a lower bound on $W_{m,r}$ for small values of r , but we do not need it here.

We now show how Kaufman et al. deduced Theorem 3.1 from Lemma 3.2. The first idea is to set $\delta = 2^{-r-1}$. The point is that there is at most one degree- r polynomial f at distance δ from the function $\mathcal{A}(x; Y_1, \dots, Y_t, \Delta_{Y_1}f(\cdot), \dots, \Delta_{Y_t}f(\cdot))$. Indeed, by the triangle inequality, the distance between any two polynomials that are δ -close to $\mathcal{A}(x; Y_1, \dots, Y_t, \Delta_{Y_1}f(\cdot), \dots, \Delta_{Y_t}f(\cdot))$ is at most $2\delta < 2^{-r}$, which is smaller than the minimum distance of $RM(m, r)$. Hence, to bound the number of polynomials $f \in RM(m, r)$ of relative weight at most $\text{wt}(f) \leq (1 - \varepsilon)2^{-\ell}$, it is enough to bound the possible number of functions of the form $\mathcal{A}(x; Y_1, \dots, Y_t, \Delta_{Y_1}f(\cdot), \dots, \Delta_{Y_t}f(\cdot))$ for the appropriate t .

The second step in the proof of Kaufmann et al. is to give an upper bound on the number of expressions of the form $\mathcal{A}(x; Y_1, \dots, Y_t, \Delta_{Y_1}f(\cdot), \dots, \Delta_{Y_t}f(\cdot))$. Since \mathcal{A} is fixed, they only have to bound the number of sets Y_i and the number of polynomials of the form $\Delta_{Y_i}f$ and raise it to the power t . They now use the fact that $\Delta_{Y_i}f$ is a polynomial of degree at most $r - \ell$ so the number of such polynomials is $2^{\binom{m}{\leq r-\ell}}$.

Combining everything, and letting $\delta = 2^{-r-2}$ so that

$$t = O(r \log(1/\varepsilon) + r^2)$$

and

$$W_{m,r}((1 - \varepsilon)2^{-\ell}) \leq \left(2^{m\ell} \cdot 2^{\binom{m}{\leq r-\ell}}\right)^t = \left(2^{m\ell + \binom{m}{\leq r-\ell}}\right)^{O(r \log(1/\varepsilon) + r^2)}. \quad (2)$$

One downside of the result of [KLP12] is that due to the dependence on r of the constant in the big O , their estimate is tight only for constant r , and becomes trivial at $r = \tilde{O}(\sqrt{m})$. Indeed, the bound in the exponent goes down roughly like $(r/m)^\ell$. Hence, the maximum is obtained for small values of ℓ , i.e., $\ell = 1$ or $\ell = 2$. For these values, the term $\log(1/\varepsilon)r + r^2$ in the exponent basically eliminates any saving that comes from $(r/m)^\ell$. Thus, to improve the bound on the weight distribution it is crucial to improve the bound for small values of ℓ . Our result does exactly this, we are able to replace the power of r in the exponent with a power of ℓ , which gives the required saving for small values of ℓ .

We now explain how we modify the arguments of [KLP12] in order to tighten the estimate given in Theorem 3.1, that hold for a broader range of parameters.

Our first observation is that one can relax the setting of δ . We set $\delta = (1 - \varepsilon)2^{-\ell-2}$, instead of $\delta = 2^{-r-2}$. The effect is that now there can be many polynomials g that are δ -close to $\mathcal{A}(x; Y_1, \dots, Y_t, \Delta_{Y_1}f(\cdot), \dots, \Delta_{Y_t}f(\cdot))$. Indeed, all we know is that the distance between any two such polynomials is at most $2\delta = (1 - \varepsilon)2^{-\ell-1}$. The point is that the number of such polynomials can be bounded from above by $W_{m,r}(2^{-\ell-1})$ which is relatively small compared to $W_{m,r}(2^{-\ell})$ and so we can (almost) think of it as 1. The effect on the expression (2) is that in the expression for t , we can (almost) replace r by ℓ . As explained before, this gives a significant saving over the bound of [KLP12].

Our second improvement comes from the simple observation that $\Delta_{Y_i}f$ can be defined by its value on the quotient space $\mathbb{F}_2^m / \langle Y_i \rangle$. As this is a space of dimension $m - \ell$, for a fixed Y_i , we can upper bound the number of polynomials of the form $\Delta_{Y_i}f$ by $2^{\binom{m-\ell}{\leq r-\ell}}$, instead of $2^{\binom{m}{\leq r-\ell}}$, which again yields a tighter estimate.

We now state our bound on the weight distribution of Reed-Muller codes.

Theorem 3.3. *Let $1 \leq \ell \leq r - 1$ and $0 < \varepsilon \leq 1/2$. Then, if $r \leq m/4$,*

$$W_{m,r}((1 - \varepsilon)2^{-\ell}) \leq (1/\varepsilon)^{8c\ell^4 \binom{m-\ell}{\leq r-\ell}},$$

where c is an absolute constant (same as in Lemma 3.2).

Proof. We shall prove by induction a stronger statement, namely,

$$W_{m,r}((1-\varepsilon)2^{-\ell}) \leq (1/\varepsilon)^{2c(m(r+3)^3(r-\ell)+(\ell+3)^2\binom{m-\ell}{r-\ell})}.$$

Set $\delta = (1-\varepsilon)2^{-\ell-2}$ and $t = c \cdot (\log(1/\delta) \log(1/\varepsilon) + \log^2(1/\delta)) \leq c \cdot \log(1/\varepsilon) \cdot (\ell+3)^2$. By Lemma 3.2, for any $f \in RM(m, r)$ of weight $\text{wt}(f) \leq (1-\varepsilon)2^{-\ell}$, there is a choice of sets $Y_1, \dots, Y_t \in (\mathbb{F}_2^m)^\ell$ such that

$$\Pr_{x \in \mathbb{F}_2^m} [\mathcal{A}(x, (Y_1, \dots, Y_t); \Delta_{Y_1} f(\cdot), \dots, \Delta_{Y_t} f(\cdot)) = f(x)] \geq 1 - \delta.$$

We next bound the number of functions g of the form

$$g = \mathcal{A}(x, (Y_1, \dots, Y_t); \Delta_{Y_1} f(\cdot), \dots, \Delta_{Y_t} f(\cdot)).$$

We can upper bound the number of sets $Y \in (\mathbb{F}_2^m)^\ell$ with $2^{m\ell}$. For each such Y , since $\Delta_Y f$ is a polynomial of degree $\leq r - \ell$ that is defined by its values on the space $\mathbb{F}_2^m / \langle Y_i \rangle$, there are at most $2^{\binom{m-\ell}{\leq r-\ell}}$ polynomials of the form $\Delta_Y f$. Thus, the number of possible such functions g is at most

$$\left(2^{m\ell} 2^{\binom{m-\ell}{\leq r-\ell}}\right)^t = (1/\varepsilon)^{c(\ell+3)^2\left(m\ell + \binom{m-\ell}{\leq r-\ell}\right)}.$$

Given any such g , the number of polynomials $g' \in RM(m, r)$ at distance at most $(1-\varepsilon)2^{-\ell-2}$ from g is at most $W_{m,r}((1-\varepsilon)2^{-\ell-1})$. Indeed, fix some f close to g . Then any other such polynomial g' has distance at most $2(1-\varepsilon)2^{-\ell-2}$ from f , and so $\text{wt}(f - g') \leq (1-\varepsilon)2^{-\ell-1}$ and $f - g' \in RM(m, r)$. Concluding we get

$$W_{m,r}((1-\varepsilon)2^{-\ell}) \leq (1/\varepsilon)^{c(\ell+3)^2\left(m\ell + \binom{m-\ell}{\leq r-\ell}\right)} \cdot W_{m,r}((1-\varepsilon)2^{-\ell-1}).$$

Since $W_{m,r}((1-\varepsilon)2^{-r}) = 1$ (as only the 0 polynomial has such small weight) and $\binom{m-\ell}{\leq r-\ell} \leq \binom{m}{\leq r} \cdot \left(\frac{r}{m}\right)^\ell$, we get by induction that

$$\begin{aligned} W_{m,r}((1-\varepsilon)2^{-\ell}) &\leq (1/\varepsilon)^{c(\ell+3)^2\left(m\ell + \binom{m-\ell}{\leq r-\ell}\right)} \cdot W_{m,r}((1-\varepsilon)2^{-\ell-1}) \\ &\leq (1/\varepsilon)^{c(\ell+3)^2\left(m\ell + \binom{m-\ell}{\leq r-\ell}\right)} \cdot (1/\varepsilon)^{2c\left(m(r+3)^3(r-\ell-1) + (\ell+4)^2\binom{m-(\ell+1)}{\leq r-(\ell+1)}\right)} \\ &\leq (1/\varepsilon)^{c(\ell+3)^2\left(m\ell + \binom{m-\ell}{\leq r-\ell}\right)} \cdot (1/\varepsilon)^{2c\left(m(r+3)^3(r-\ell-1) + (\ell+4)^2\binom{m-\ell}{\leq r-\ell} \cdot \left(\frac{r}{m}\right)\right)} \\ &\leq (1/\varepsilon)^{2cm\left((r+3)^3(r-\ell)\right)} \cdot (1/\varepsilon)^{c\binom{m-\ell}{r-\ell}\left((\ell+3)^2 + 2(\ell+4)^2\frac{r}{m}\right)} \\ &\leq^* (1/\varepsilon)^{2c\left(m(r+3)^3(r-\ell) + (\ell+3)^2\binom{m-\ell}{r-\ell}\right)}, \end{aligned}$$

where in inequality (*) we use the fact that $r < m/4$. The bound in the statement of the theorem follows by a simple manipulation. This concludes the proof of the theorem. \square

4 Random submatrices of $E(m, r)$

As discussed in the introduction and Section 2, in order to understand the ability to decode from erasures it is important to understand the following questions. Consider randomly chosen set U of a given parameter size k :

Question 4.1. *What is the largest s for which the submatrix U^r has full column-rank with high probability?*

Question 4.2. *What is the smallest s for which the submatrix U^r has full row-rank with high probability?*

In this section we provide an answer to each of these questions.¹⁶ Note that for any degree- r , the number of rows of $E(m, r)$, namely $\binom{m}{\leq r}$, is an upper bound on the value of s for the first question and a lower bound for the second. For small r we prove that we can approach this optimal bound asymptotically in both.

Note that, interestingly, the duality property of RM codes allows to relate question 4.1 and 4.2 to each other but for different ranges of the parameters. Namely, the following holds.

Lemma 4.3. *For a set $S \subseteq [n]$, denote by $E(m, d)[S]$ the sub-matrix of $E(m, d)$ obtained by selecting the columns indexed by S . For any $s \leq n$,*

$$\left\{ S \in \binom{[n]}{s} : \text{rk}(E(m, d)[S]) = s \right\} = \left\{ S \in \binom{[n]}{s} : \text{rk}(E(m, m-d-1)[S^c]) = n - \binom{m}{\leq d} \right\}.$$

Note that $E(m, d)[S] = s$ means that $E(m, d)[S]$ has full column-rank and $E(m, m-d-1)[S^c] = n - \binom{m}{\leq d}$ means that $E(m, m-d-1)[S^c]$ has full row-rank.

Corollary 4.4. *For an integer $s \in [n]$, denote by $E(m, d)[s]$ the random matrix obtained by sampling s columns uniformly at random in $E(m, d)$. Then,*

$$\Pr\{\text{rk}(E(m, d)[s]) = s\} = \Pr\left\{\text{rk}(E(m, m-d-1)[n-s]) = n - \binom{m}{\leq d}\right\},$$

where both terms are the probability of drawing a uniform erasure pattern of size s which can be corrected with the code $\ker E(m, d)$.

This correspondence follows from Lemmas 2.8, 2.6 and 2.11. We provide the proof below for convenience.

Proof of Lemma 4.3. Note that

$$\begin{aligned} & \{S \in \binom{[n]}{s} : \text{rk}(E(m, d)[S]) < s\} \\ & \equiv \{S \in \binom{[n]}{s} : \exists z \in \ker(E(m, d)), \text{ s.t. } \text{supp}(z) \subseteq S, z \neq 0\} \\ & \equiv \{S \in \binom{[n]}{s} : \exists z \in \ker(E(m, d)) \text{ s.t. } z[S^c] = 0, z \neq 0\}, \end{aligned}$$

and using Lemma 2.11, previous set is equal to

$$\begin{aligned} & \{S \in \binom{[n]}{s} : \exists z \in \text{Im}(E(m, m-d-1)) \text{ s.t. } z[S^c] = 0, z \neq 0\} \\ & \equiv \{S \in \binom{[n]}{s} : E(m, m-d-1)[S^c] \text{ is not full row-rank}\} \\ & \equiv \{S \in \binom{[n]}{s} : \text{rk}(E(m, m-d-1)[n-s]) < n - \binom{m}{\leq d}\}. \end{aligned}$$

□

This equivalence property implies that it is sufficient to answer each question in one of the two extremal regimes, which we next cover.

¹⁶Using tensoring to produce linearly independent vectors has also been studied recently in the context of real vectors [BCMV14].

4.1 Random submatrices of $E(m, r)$, for small r , have full column-rank

The following theorem addresses Question 4.1 in the case of low degree- r .

Theorem 4.5. *Let $\varepsilon > 0$ and k, m, r integers such that $s < \binom{m - \log(\binom{m}{\leq r}) - \log(1/\varepsilon)}{\leq r}$. Then, with probability larger than $1 - \varepsilon$ if we pick $u_1, \dots, u_s \in \mathbb{F}_2^m$ uniformly at random we get that the evaluation vectors, u_1^r, \dots, u_s^r are linearly independent.*

Observe that for $r = o(\sqrt{m/\log m})$ the bound on s is $(1 - o(1))\binom{m}{\leq r}$, which will give us a capacity-achieving result.

As discussed in Section 1.4 (Theorem 1.3), to prove the theorem we have to understand the set of common zeroes of degree- r polynomials. More accurately, we need to give an upper bound on the number of common zeroes of polynomials in some linear space.

We start by introducing some notation and then discuss the reduction from Theorem 4.5 to the problem of determining the number of common zeroes of a space of polynomials.

Given a set of points $u_1, \dots, u_s \in \mathbb{F}_2^m$ we define

$$\mathcal{I}(u_1, \dots, u_s) = \{f \in \mathbb{P}(m, r) \mid \forall i \ f(u_i) = 0\}.$$

When U is an $m \times s$ matrix we define $\mathcal{I}(U) = \mathcal{I}(u_1, \dots, u_s)$, where u_i is the i th column of U . It is clear that $\mathcal{I}(U)$ is a vector space. Similarly, for a set of polynomials $F \subseteq \mathbb{P}(m, r)$ we denote

$$\mathcal{V}(F) = \{u \in \mathbb{F}_2^m \mid \forall f \in F \ f(u) = 0\}.$$

In other words, $\mathcal{V}(F)$ is the set of common zeroes of F . From the definition it is clear that if $F_1 \subseteq F_2$ then $\mathcal{V}(F_2) \subseteq \mathcal{V}(F_1)$ and similarly, if $U_1 \subseteq U_2$ then $\mathcal{I}(U_2) \subseteq \mathcal{I}(U_1)$.

The next lemmas explore the connection between the dual space of U^r , $\mathcal{I}(U)$ and $\mathcal{V}(\mathcal{I}(U))$. Hereafter we interpret a vector f of length $\binom{m}{\leq r}$ as a polynomial in $\mathbb{P}(m, r)$, by viewing its coordinates as coefficients of the relevant monomials. We abuse notation and call this polynomial f as well.

Lemma 4.6. *Let U be an $m \times s$ matrix. Then, a vector f of length $\binom{m}{\leq r}$ satisfies $f \cdot U^r = 0$ if and only if the corresponding polynomial $f(x_1, \dots, x_m)$ is in $\mathcal{I}(U)$, namely, $f \in \mathcal{I}(U)$.*

Proof. The proof is immediate from the correspondence between vectors to polynomials and from the definition of U^r . Indeed, for a column u_i we have that the coordinates of u_i^r correspond to all evaluations of monomials of degree $\leq r$ on u_i . Similarly, the coordinates of the vector f correspond to coefficients of the polynomial $f(x_1, \dots, x_m)$. Thus, $f \cdot u_i^r$ is equal to $f(u_i)$. Hence, $f \cdot U^r = 0$ if and only if $f(u_1) = \dots, f(u_s) = 0$, i.e. if and only if $f \in \mathcal{I}(U)$. \square

Lemma 4.7. *Let U be an $m \times s$ binary matrix. Then, for any $u \in \mathbb{F}_2^m$ we have that u^r is in the linear span of the columns of U^r if and only if*

$$\mathcal{I}(U) = \mathcal{I}(U \cup \{u\}),$$

namely, every degree $\leq r$ polynomial that vanishes on the columns of U also vanishes on u .

Proof. It is clear that u^r linearly depends on the columns of U^r if and only if for every vector f such that $f \cdot U^r = 0$, it holds that $f \cdot u^r = 0$, namely, that $f(u) = 0$. By Lemma 4.6 this is equivalent to saying that $\mathcal{I}(U) = \mathcal{I}(U \cup \{u\})$. \square

Similarly, we get an equivalence when consider the common zeros of the polynomials that vanish on the columns of U .

Lemma 4.8. *We have that $u \in \mathcal{V}(\mathcal{I}(U))$ if and only if u^r is spanned by the columns of U^r .*

Proof. If u^r is spanned by the columns of U^r then by Lemma 4.7 $\mathcal{I}(U) = \mathcal{I}(U \cup \{u\})$. Thus, $u \in \mathcal{V}(\mathcal{I}(U \cup \{u\})) = \mathcal{V}(\mathcal{I}(U))$. Conversely, if $u \in \mathcal{V}(\mathcal{I}(U))$ then $\mathcal{I}(U) \subseteq \mathcal{I}(U \cup \{u\})$. As $\mathcal{I}(U \cup \{u\}) \subseteq \mathcal{I}(U)$ we get $\mathcal{I}(U \cup \{u\}) = \mathcal{I}(U)$ and by Lemma 4.7 it follows that u^r is spanned by the columns of U^r . \square

Finally, we make the following simple observation.

Lemma 4.9. $\mathcal{I}(\mathcal{V}(\mathcal{I}(U))) = \mathcal{I}(U)$.

Proof. Denote $\mathcal{V} = \mathcal{V}(\mathcal{I}(U))$. It is clear that $U \subseteq \mathcal{V}$ and hence $\mathcal{I}(\mathcal{V}) \subseteq \mathcal{I}(U)$. On the other hand, let $f \in \mathcal{I}(U)$ and $v \in \mathcal{V}$. Lemmas 4.7 and 4.8 imply that $\mathcal{I}(U) = \mathcal{I}(U \cup \{v\})$. Thus, $f(v) = 0$ and hence $\mathcal{I}(U) \subseteq \mathcal{I}(\mathcal{V})$. \square

Going back to our original problem, assume that we picked s columns at random and got linearly independent evaluation vectors. Now we have to understand the probability that a randomly chosen $u \in \mathbb{F}_2^m$ will give an independent evaluation vector. By Lemma 4.8, this amounts to understanding the probability that u belongs to $\mathcal{V} := \mathcal{V}(\mathcal{I}(U))$. By linear algebra arguments we have the following identity

$$\dim(\mathcal{I}(U)) = \binom{m}{\leq r} - \text{rank}(U^r).$$

Thus, our goal is understanding how many common zeroes can the polynomials in an $\left(\binom{m}{\leq r} - s\right)$ -dimensional space have.

The way to prove that a given set of polynomials does not have too many common zeros is to show that any large set of points (in our case, \mathcal{V}) has many linearly independent degree- r polynomials that are defined over it. That is, we only consider the restriction of polynomials of degree- r to the points in \mathcal{V} and we identify two polynomials if they are equal when restricted to \mathcal{V} . Notice that this is the same as showing that the rank of $E(m, r)[\mathcal{V}]$ is large, i.e., that there are many linearly independent columns that are indexed by elements of \mathcal{V} . Thus, two such polynomials f, g are identified if and only if $f - g \in \mathcal{I}(\mathcal{V})$. Stated differently, we wish to show that the dimension of the quotient space $\mathbb{P}(m, r)/\mathcal{I}(\mathcal{V}) = \mathbb{P}(m, r)/\mathcal{I}(U)$ (by Lemma 4.9) is large. Indeed, if we can lower bound this dimension in terms of \mathcal{V} then, since

$$\begin{aligned} \text{rank}(E(m, r)[\mathcal{V}]) &= \dim(\mathbb{P}(m, r)/\mathcal{I}(\mathcal{V})) = \dim(\mathbb{P}(m, r)/\mathcal{I}(U)) = \dim(\mathbb{P}(m, r)) - \dim(\mathcal{I}(U)) \\ &= \binom{m}{\leq r} - \left(\binom{m}{\leq r} - s\right) = s, \end{aligned}$$

we will get that some function of $|\mathcal{V}|$ is upper bounded by s . Thus, unless s is large, $|\mathcal{V}|$ is small and hence the probability that a randomly chosen u is independent of U^r is high. We state our main lemma next.

Lemma 4.10. *Let $\mathcal{V} \subseteq \mathbb{F}_2^m$ such that $|\mathcal{V}| > 2^{m-t}$. Then there are more than $\binom{m-t}{\leq r}$ linearly independent degree $\leq r$ polynomials defined on \mathcal{V} , i.e., $\dim(E(m, r)[\mathcal{V}]) > \binom{m-t}{\leq r}$.*

We give two different proofs of this fact. The first uses a hashing argument; if \mathcal{V} is large, then, after some linear transformation, its projection on a set of roughly $\log(|\mathcal{V}|)$ many coordinates is full. Thus, it supports at least $\binom{\log(|\mathcal{V}|)}{\leq r}$ many linearly independent degree- r monomials. The second proof relies on a somewhat tighter bound that was obtained by Wei [Wei91], who studied the *generalized Hamming weight* of Reed-Muller codes. As Wei's result gives slightly tighter bounds compared to the hashing argument (although both lead to a capacity-achieving result), this is the

proof that we give in the main body of the paper. For completeness, and as the hashing argument is more self contained we give it in Appendix B.

We start by discussing the notion of generalized Hamming weight. Let $C \subseteq \mathbb{F}_2^n$ be a linear code and $D \subseteq C$ a linear subcode. We denote

$$\text{supp}(D) = \{i : \exists y \in D, \text{ such that } y_i \neq 0\}.$$

In other words, the support of D is the union of the supports of all codewords in D .

Definition 4.11 (Generalized Hamming weight). *For a code C of length n and an integer a we define*

$$d_a(C) = \min\{\text{supp}(D) \mid D \subseteq C \text{ is a linear subcode with } \dim(D) = a\}.$$

Thus, $d_a(C)$ is the minimal size of a set of coordinates S , such that there exists a subcode D , of dimension $\dim(D) = a$, that is supported on S . The reason for this definition is that for any code C if we let $a = \dim(C)$ then $d_a(C) = n - d$, where d is the minimal distance of C . By considering the complement set S^c the next lemma gives an equivalent definition of $d_a(C)$.

Lemma 4.12. *For a code C of length n and an integer a we have that*

$$d_a(C) = \max\{b \mid \forall |S| < b \text{ we have that } \dim(C[S^c]) > \dim(C) - a\}.$$

Proof. The proof follows immediately from a simple linear algebra argument. If D is a subcode of C that is supported on a set of coordinates S then $\dim(C) = \dim(D) + \dim(C[S^c])$. \square

The alternative definition given in Lemma 4.12 is very close to what we need. We wish to show that for any large \mathcal{V} , there are many linearly independent degree $\leq r$ polynomials that are defined on \mathcal{V} . In other words, we wish to prove that

$$d_{o(1)\binom{m}{\leq r}}(RM(m, r)) \geq 2^m - \varepsilon \cdot 2^m / \binom{m}{\leq r}.$$

Indeed, this will imply that for any $|\mathcal{V}| \geq \varepsilon \cdot 2^m / \binom{m}{\leq r}$ there are at least $(1 - o(1))\binom{m}{\leq r}$ linearly independent degree $\leq r$ monomials defined on \mathcal{V} (\mathcal{V} plays the role of S^c in Lemma 4.12).

The next theorem of Wei [Wei91] computes exactly the generalized Hamming weight of Reed-Muller codes. For stating the theorem we need the following technical claim.

Lemma 4.13 (Lemma 2 of [Wei91]). *For every $0 \leq a \leq \binom{m}{\leq r}$ there is a unique way of expressing a as $a = \sum_{i=1}^{\ell} \binom{m_i}{\leq r_i}$, where $m_i - r_i = m - r - i + 1$.*

Theorem 4.14 ([Wei91]). *Let $0 \leq a \leq \binom{m}{\leq r}$ be an integer. Then, $d_a(RM(m, r)) = \sum_{i=1}^{\ell} 2^{m_i}$, where $a = \sum_{i=1}^{\ell} \binom{m_i}{\leq r_i}$ is the unique representation of a according to Lemma 4.13.*

We are now ready to prove Lemma 4.10.

Proof of Lemma 4.10. For $a = \sum_{i=1}^t \binom{m-i}{\leq r-1}$, Theorem 4.14 implies that $d_a(RM(m, r)) = \sum_{i=1}^t 2^{m-i} = 2^m - 2^{m-t}$. Thus, if $|\mathcal{V}| > 2^m - d_a(RM(m, r)) = 2^{m-t}$ then there are more than $\binom{m}{\leq r} - a = \binom{m}{\leq r} - \sum_{i=1}^t \binom{m-i}{\leq r-1}$ many linearly independent degree- r polynomials defined on \mathcal{V} . To make sense of parameters we shall need the following simple calculation. We give the straightforward proof in Section A.

Claim 4.15. $\binom{m}{\leq r} - \sum_{i=1}^t \binom{m-i}{\leq r-1} = \binom{m-t}{\leq r}$.

Thus, if $|\mathcal{V}| > 2^{m-t}$ then there are more than $\binom{m-t}{\leq r}$ linearly independent degree $\leq r$ polynomials defined on \mathcal{V} . \square

We can now bound the number of common zeroes of all polynomials that vanish on a given set U .

Lemma 4.16. *Let $\varepsilon > 0$ be a constant and s an integer such that $s < \binom{m - \lceil \log(\binom{m}{\leq r}) + \log(1/\varepsilon) \rceil}{\leq r}$. Let U be an $m \times s$ matrix such that $\text{rank}(U^r) = s$, namely, the columns of U^r are linearly independent. Then, $|\mathcal{V}(\mathcal{I}(U))| \leq \varepsilon \cdot 2^m / \binom{m}{\leq r}$.*

Proof. By the discussion above we know that $\dim(\mathbb{P}(m, r)/\mathcal{I}(U)) = s$. I.e., the dimension of the space of degree $\leq r$ polynomials that are defined on $\mathcal{V} = \mathcal{V}(\mathcal{I}(U))$ is s .

Let t be the minimal integer such that $2^{m-t} \leq \varepsilon \cdot \frac{2^m}{\binom{m}{\leq r}}$, i.e., $t = \lceil \log\left(\binom{m}{\leq r}\right) + \log(1/\varepsilon) \rceil$. Assume towards a contradiction that $|\mathcal{V}| > \varepsilon \cdot \frac{2^m}{\binom{m}{\leq r}} \geq 2^{m-t}$. Lemma 4.10 implies that there are more than $\binom{m-t}{\leq r}$ linearly independent polynomials defined on \mathcal{V} .

As there are at most s polynomial defined on \mathcal{V} we must have

$$s > \binom{m-t}{\leq r} = \binom{m - \lceil \log(\binom{m}{\leq r}) + \log(1/\varepsilon) \rceil}{\leq r},$$

in contradiction to the assumption of the lemma. \square

We can now derive Theorem 4.5, the main result of this subsection.

Proof of Theorem 4.5. We start by picking random points one after the other. Assume that we picked $\ell < s$ points and they gave rise to linearly independent evaluation vectors. By Lemmas 4.8 and 4.16, as long as $s \leq \binom{m - \lceil \log(\binom{m}{\leq r}) + \log(1/\varepsilon) \rceil}{\leq r}$, the probability that the next random point will not generate an independent evaluation vector is at most $\varepsilon / \binom{m}{\leq r}$. Repeating this argument $\binom{m - \lceil \log(\binom{m}{\leq r}) + \log(1/\varepsilon) \rceil}{\leq r} < \binom{m}{\leq r}$ times, the probability that we do not get a set of independent evaluation vectors is at most ε . \square

4.2 Random submatrices of $E(m, r)$, for small r , have full row-rank

In this section we study Question 4.2 for the case of low degree. Thus, our goal is proving that, with high probability, a random set U of columns of $E(m, r)$, of the appropriate size, has full row-rank. As we showed in Corollary 4.4, this is equivalent to studying Question 4.1 in the case of high degree. We prove the following theorem.

Theorem 4.17. *Let $0 < \delta < 1/3$, $\eta = O(1/\log(1/\delta))$ and m, r integers such that $r \leq \eta m$. Let $s = \lceil (1 + \delta) \binom{m}{\leq r} \rceil$. Then, except of probability not larger than $\exp\left(-\Omega\left(\min(\delta, 2^{-r}) \cdot \binom{m}{\leq r}\right)\right)$, if we pick $u_1, \dots, u_s \in \mathbb{F}_2^m$ uniformly at random we get that the evaluation vectors, u_1^r, \dots, u_s^r span the entire vector space $\mathbb{F}_2^{\binom{m}{\leq r}}$.*

As in the proof of Theorem 4.5 we will consider the dual space to evaluation vectors - the space of degree- r polynomials. Our proof strategy is to show that for every possible polynomial of degree- r , with high probability, there is at least one point in our set on which the polynomial does

not vanish. Thus, by the discussion in Section 4.1, this means that the dual space contains only the zero polynomial, which is exactly what we wish to prove. To apply this strategy we will need to know the number of polynomials that have a certain number of nonzeros. Such an estimate was given in Theorem 3.3, following [KLP12].

Proof. Set $\varepsilon = \delta/4$. For an integer $1 \leq \ell \leq r$ we denote with P_ℓ the set of degree $\leq r$ polynomials whose fraction of nonzeros is between $(1 - \varepsilon)2^{-\ell-1}$ and $(1 - \varepsilon)2^{-\ell}$. By Theorem 3.3,

$$|P_\ell| \leq W_{m,r}((1 - \varepsilon)2^{-\ell}) \leq (1/\varepsilon)^{8c\ell^4 \binom{m-\ell}{\leq r-\ell}}.$$

Let f be some polynomial in P_ℓ . When picking s points at random, the probability that f vanishes on all of them is at most $(1 - (1 - \varepsilon)2^{-\ell-1})^s$. Thus, the probability that any $f \in P_\ell$ vanishes on all s points is at most $|P_\ell| \cdot (1 - (1 - \varepsilon)2^{-\ell-1})^s$. We would like to show that this probability is small, when ranging over all ℓ . We first study the case that $\ell > 0$, i.e., of the contribution from the polynomials that have at most $(1 - \varepsilon)/2$ nonzeros. Although in this case the probability of hitting a nonzero gets smaller and smaller as ℓ grows, the size of P_ℓ goes down at a faster rate (by Theorem 1.5), and therefore we get an exponentially small probability of missing some polynomial. When $\ell = 0$, although the number of polynomials of such high weight is huge, the probability of missing any of them is tiny, and so we are ok in this case as well. We now do the formal calculation.

By the above, the probability that some polynomial whose fraction of nonzeros is at most $(1 - \varepsilon)/2$ vanishes on all s points is at most

$$\sum_{\ell=1}^{r-1} |P_\ell| \cdot (1 - (1 - \varepsilon)2^{-\ell-1})^s. \quad (3)$$

Indeed, any degree- r polynomial is nonzero with probability at least 2^{-r} and hence the above summation goes over all such polynomials.

Let us estimate a typical summand,

$$\begin{aligned} |P_\ell| \cdot \left(1 - (1 - \varepsilon)2^{-\ell-1}\right)^s &\leq (1/\varepsilon)^{8c\ell^4 \binom{m-\ell}{\leq r-\ell}} \cdot (1 - (1 - \varepsilon)2^{-\ell-1})^s \\ &\leq (1/\varepsilon)^{8c\ell^4 \binom{m-\ell}{\leq r-\ell}} \cdot \exp(-(1 - \varepsilon)2^{-\ell-1}s) \\ &\leq (1/\varepsilon)^{8c\ell^4 \binom{m-\ell}{\leq r-\ell}} \cdot \exp\left(-(1 + \delta)(1 - \varepsilon) \binom{m}{\leq r} 2^{-\ell-1}\right). \end{aligned}$$

Let $\eta = O(1/\log(1/\varepsilon)) = O(1/\log(1/\delta))$. From the fact that

$$\binom{m-\ell}{\leq r-\ell} \leq \binom{m}{\leq r} \cdot \left(\frac{r}{m}\right)^\ell,$$

we get by simple manipulations that if $r \leq \eta m$ then

$$\begin{aligned} |P_\ell| \cdot \left(1 - (1 - \varepsilon)2^{-\ell-1}\right)^s &\leq \exp\left(\binom{m}{\leq r} \cdot \left(3^{-\ell} - (1 + \delta)(1 - \varepsilon)2^{-\ell-1}\right)\right) \\ &\leq \exp\left(-\Omega\left(\binom{m}{\leq r} \cdot 2^{-\ell}\right)\right). \end{aligned}$$

Going back to Equation (3) we see that

$$\sum_{\ell=1}^{r-1} |P_\ell| \cdot (1 - (1 - \varepsilon)2^{-\ell-1})^s \leq \sum_{\ell=1}^{r-1} \exp\left(-\Omega\left(\binom{m}{\leq r} \cdot 2^{-\ell}\right)\right) \leq \exp\left(-\Omega\left(2^{-r} \binom{m}{\leq r}\right)\right).$$

All we have to do now is bound from above the probability that there exists a polynomial that has more than $(1 - \varepsilon)/2$ fraction of nonzeros, but that vanishes on all the s chosen points. As there are at most $2^{\binom{m}{\leq r}}$ such polynomials, this probability can be bounded from above by

$$\begin{aligned}
2^{\binom{m}{\leq r}} \cdot \left(1 - \frac{1 - \varepsilon}{2}\right)^s &= 2^{\binom{m}{\leq r}} \cdot \left(\frac{1 + \varepsilon}{2}\right)^s \\
&\leq 2^{\binom{m}{\leq r}} \cdot \left(\frac{1 + \varepsilon}{2}\right)^{(1+\delta)\binom{m}{\leq r}} \\
&= \left((1 + \varepsilon)^{1/\delta} \cdot \left(\frac{1 + \varepsilon}{2}\right)\right)^{\delta\binom{m}{\leq r}} \\
&\stackrel{(*)}{\leq} \left(e^{1/4} \cdot \left(\frac{1 + \varepsilon}{2}\right)\right)^{\delta\binom{m}{\leq r}} \\
&\stackrel{(\dagger)}{\leq} \exp\left(-\frac{1}{3}\delta\binom{m}{\leq r}\right),
\end{aligned}$$

where inequality $(*)$ holds since $\varepsilon = \delta/4$ and inequality (\dagger) follows by our choice $\delta < 1/3$. Concluding, when picking s points at random, the probability that there will be some polynomial of degree $\leq r$ that does not vanish on any of the chosen points is at most

$$\exp\left(-\Omega\left(2^{-r}\binom{m}{\leq r}\right)\right) + \exp\left(-\frac{1}{3}\delta\binom{m}{\leq r}\right) = \exp\left(-\Omega\left(\min(\delta, 2^{-r})\binom{m}{\leq r}\right)\right).$$

By simple linear algebra and Lemma 4.6 it follows that if no such polynomial vanishes on all points that we picked then their evaluation vectors span the entire space. This concludes the proof of the theorem. \square

5 Reed-Muller code for erasures

5.1 Low-rate regime

Recall that from Corollary 2.7, if G is an $n \times k$ generator matrix of a code, then the code can correct s random erasures if erasing a random subset of $n - s$ rows in G (which gives the random matrix $G_{n-s, \cdot}$) has full span, i.e.,

$$\Pr\{\text{rank}(G_{n-s, \cdot}) = k\} \rightarrow 1, \quad \text{as } n \rightarrow \infty.$$

Hence, since the transpose of $E(m, r)$ is an $n \times k$ generator matrix for $RM(m, r)$, we get our result for low-rate RM codes on the BEC as a direct consequence of Theorem 4.17.

Corollary 5.1. *Let $0 < \delta < 1/3$, $\eta = O(1/\log(1/\delta))$ and m, r integers such that $r \leq \eta m$. Then, $RM(m, r)$ can correct $2^m - (1 + \delta)\binom{m}{\leq r}$ random erasures, i.e., it is δ -close to capacity-achieving. Moreover, if $r = o(m)$, then $RM(m, r)$ is capacity-achieving on the BEC.*

5.2 High-rate regime

We now use the parity-check matrix interpretation of correcting errors, namely from Corollary 2.9, a code with parity-check matrix H can correct s errors if the random set of s columns (which we denote $H[s]$) are linearly independent with high probability, i.e.,

$$\Pr\{\text{rk}(H[s]) = s\} \rightarrow 1, \quad \text{as } n \rightarrow \infty.$$

Using now Theorem 4.5 and the fact that $E(m, r)$ is a parity-check matrix for $RM(m, m - r - 1)$ (see Lemma 2.11), we obtain our result for the performance of high-rate RM codes on the BEC.

Corollary 5.2. *Let $\varepsilon > 0$, $r \leq m$ be two positive integers and $s = \lfloor \binom{m - \log(\binom{m}{\leq r}) - \log(1/\varepsilon)}{\leq r} \rfloor$. Then, $RM(m, m - r)$ can correct s random erasures with probability larger than $1 - \varepsilon$. In particular, if $m - r = o(\sqrt{m/\log m})$, then $RM(m, r)$ is capacity-achieving on the BEC.*

The following calculation gives a better sense of the parameters (the proof is in Section A).

Claim 5.3. *For $r < \sqrt{\frac{\delta m}{4 \log(m)}}$ and $\varepsilon > m^{-r/2}$ we have that $\binom{m - \log(\binom{m}{\leq r}) - \log(1/\varepsilon)}{\leq r} > (1 - \delta) \binom{m}{\leq r}$.*

6 Reed-Muller code for errors

We present next results for errors at both low and high rate. The results at low-rate rely on the weight distribution results of Section 3, whereas the high-rate results rely on a novel relation between decoding from errors and decoding from erasures.

6.1 Low-rate regime

Theorem 6.1. *Let $\delta > 0$. There exists¹⁷ $\eta = O(1/\log(1/\delta))$ such that the following holds. For any two integers r and m satisfying $r/m \leq \eta$, and any p satisfying*

$$1 - h(p) = (1 + \delta)R, \quad \text{where } R = \frac{\binom{m}{\leq r}}{n},$$

$RM(m, r)$ can correct pn random errors with probability at least $\exp\left(-\Omega\left(\min(\delta, 2^{-r}) \cdot \binom{m}{\leq r}\right)\right)$. In particular, for $r = o(m)$, $RM(m, r)$ is capacity-achieving on the BSC.

Before giving the proof we make a small calculation to get a better sense of what parameters we should expect. Since R is small we can expect to correct a fraction of errors approaching $1/2$. Let us denote $p = (1 - \xi)/2$. We now wish to figure how small should ξ be. At corruption rate close to $1/2$ we have that

$$h(p) = h(1/2 - \xi/2) = 1 - \xi^2/(2 \ln(2)) + \Theta(\xi^4). \quad (4)$$

Thus, if we wish to have $(1 + \delta)R = 1 - h(p)$ then ξ should satisfy

$$(1 + \delta)R = 1 - h(p) = \xi^2/(2 \ln(2)) - \Theta(\xi^4). \quad (5)$$

Hence, $\xi^2 = \Theta(R)$.

We now give the proof following the outline described in Section 1.4.

Proof. Let s be the number of errors, i.e., $s = pn$ and $p = 1/2 - \xi/2$. A bad error pattern $z \in \mathbb{F}_2^n$ is one for which there exists another error pattern $z' \in \mathbb{F}_2^n$, of weight s , such that $z + z'$ is a codeword in $RM(m, r)$. We concentrate on the case that $w(z') = s$ as this is the most interesting case.

Note that since both z and z' are different and have the same weight, the weight of $z + z'$ must be even and in $\{d, \dots, 2s\}$. As both $z + z'$ and the all 1 vector are codewords, we also have that the weight of $z + z'$ is at most $n - d$, hence $w(z + z') \in \{d, \dots, n - d\}$. Therefore, counting the number

¹⁷The exact value of η is given in (11).

of bad error patterns is equivalent to counting the number of weight s vectors that can be obtained by “splitting” codewords of even weight in $\{d, \dots, n-d\}$. Note that for a codeword y of weight $w(y) = w$, there are $w/2$ choices for the support of z inside $\text{supp}(y)$ and $s - w/2$ choices outside the codeword’s support. Indeed, z and z' must cancel each other outside the support of y and hence they have the same weight inside $\text{supp}(y)$.¹⁸ It follows, that for a fixed y there are

$$\binom{w}{w/2} \binom{n-w}{s-w/2}$$

possibilities to pick a bad error pattern with intersection $w/2$ with $\text{supp}(y)$. Denoting by \mathcal{B} the set of bad error patterns and $N_{m,r}(w)$ the number of codewords of weight w in $RM(m, r)$, a union bound gives

$$\Pr\{\mathcal{B}\} \leq \sum_{w \in \{d, \dots, n-d\}} N_{m,r}(w) \frac{\binom{w}{w/2} \binom{n-w}{s-w/2}}{\binom{n}{s}}.$$

We are now going to prove that $\Pr[\mathcal{B}]$ is exponentially small, for our setting of parameters, which will imply the theorem. Since for¹⁹ $\alpha \in (0, 1)$,

$$2^{nh(\alpha) - O(\log(n))} \leq \binom{n}{\alpha n} \leq 2^{nh(\alpha)},$$

and recalling the entropy approximation of (4), we have, by defining $\beta = w/n$,

$$\begin{aligned} \Pr\{\mathcal{B}\} &\leq \sum_{\beta \in \{d/n, \dots, 1-d/n\}} N_{m,r}(\beta n) \frac{2^{\beta n} 2^{n \left(1 - \beta - \frac{\xi^2}{(1-\beta)2 \ln(2)} + O\left(\frac{\xi^4}{(1-\beta)^2}\right)\right)}}{2^{n \left(1 - \frac{\xi^2}{2 \ln(2)}\right) - O(\log n)}} \\ &= \sum_{\beta \in \{d/n, \dots, 1-d/n\}} N_{m,r}(\beta n) 2^{-\frac{n\xi^2}{2 \ln(2)} \frac{\beta}{(1-\beta)} + nO\left(\frac{\xi^4}{(1-\beta)^2}\right) + O(\log n)}. \end{aligned}$$

Let $\varepsilon = \delta/3$. We next upper bound the above summation by grouping codewords of weights between $(1 - \varepsilon)2^{-\ell-1}$ and $(1 - \varepsilon)2^{-\ell}$, with $\ell \in \{1, 2, \dots, r-1\}$. For codewords of weight close to $1/2$, we use the fact that there are $2^{\binom{m}{\leq r}}$ codewords in $RM(m, r)$. Since the function $\beta \rightarrow \beta/(1 - \beta)$ is increasing, we obtain the following bound where $|P_\ell|$ is, as before, the number of codewords having weights between $(1 - \varepsilon)2^{-\ell-1}$ and $(1 - \varepsilon)2^{-\ell}$:

$$\Pr\{\mathcal{B}\} \leq \sum_{\ell \in \{1, 2, \dots, r-1\}} |P_\ell| 2^{-\frac{n\xi^2}{2 \ln(2)} \frac{(1-\varepsilon)2^{-\ell-1}}{(1-(1-\varepsilon)2^{-\ell-1})} + nO\left(\frac{\xi^4}{(1-(1-\varepsilon)2^{-\ell-1})^2}\right) + O(\log n)} \quad (6)$$

$$+ 2^{\binom{m}{\leq r}} 2^{-\frac{n\xi^2}{2 \ln(2)} \frac{(1-\varepsilon)/2}{(1-(1-\varepsilon)/2)} + nO(\xi^4/(1+\varepsilon)^2) + O(\log n)}. \quad (7)$$

Using the inequality of Theorem 3.3:

$$|P_\ell| \leq (1/\varepsilon)^{8c\ell^4 \binom{m-\ell}{\leq r-\ell}},$$

¹⁸When z and z' do not have the same weight they still have to cancel each other outside y . Thus, $\text{supp}(z)$ must have intersection $w/2 + (w(z) - w(z'))/2$ with $\text{supp}(y)$.

¹⁹The binomial coefficient should be defined for the rounding of αn with either ceiling or floor functions.

and noting that by (5), $\frac{\xi(n)^2}{2\ln(2)} = (1+\delta)R + O(\xi^4) = (1+\delta)R + \Theta(R^2)$, the term in (7) is bounded as

$$2^{-\binom{m}{\leq r} \left((1+\delta) \frac{(1-\varepsilon)/2}{1-(1-\varepsilon)/2} - 1 + \Theta(R) \right) + O(m)} = 2^{-\binom{m}{\leq r} \left((1+\delta) \frac{1-\varepsilon}{1+\varepsilon} - 1 + \Theta(R) \right) + O(m)},$$

which vanishes since $\varepsilon < \delta/2$ (recalling that $R = o(1)$) and is overall $2^{-\Omega\left(\delta\binom{m}{\leq r}\right)}$.

For the summands in (6), they are bounded as

$$(1/\varepsilon)^{8c\ell^4\binom{m-\ell}{\leq r-\ell}} 2^{-(1+\delta)\binom{m}{\leq r} \frac{(1-\varepsilon)2^{-\ell-1}}{1-(1-\varepsilon)2^{-\ell-1}} + \Theta(R) + nO\left(\frac{\xi^4}{(1-(1-\varepsilon)2^{-\ell-1})^2}\right) + O(m)}$$

where the above exponent is upper bounded

$$-\binom{m}{\leq r} \left((1+\delta) \frac{(1-\varepsilon)2^{-\ell-1}}{1-(1-\varepsilon)2^{-\ell-1}} + \Theta(R) - 8c\log(1/\varepsilon)\ell^4 \cdot \left(\frac{r}{m}\right)^\ell \right) + O(m). \quad (8)$$

Note that the second bracket in this exponent is given by

$$(1+\delta) \frac{(1-\varepsilon)2^{-\ell-1}}{(1-(1-\varepsilon)2^{-\ell-1})} - 8c\log(1/\varepsilon)\ell^4 \cdot \left(\frac{r}{m}\right)^\ell \quad (9)$$

$$= (1-\varepsilon)2^{-\ell-1} \left(\frac{1+\delta}{1-(1-\varepsilon)2^{-\ell-1}} - 16c \frac{\log(1/\varepsilon)}{1-\varepsilon} \ell^4 \cdot \left(\frac{2r}{m}\right)^\ell \right), \quad (10)$$

which is positive for $r/m < \frac{1+\delta}{8c \frac{3+\varepsilon}{1-\varepsilon} \log(1/\varepsilon)}$, and since $\varepsilon < \delta/2$, it is equal to $\Omega((1-\varepsilon)2^{-\ell-1})$ for

$$r/m \leq \eta = \frac{1+\delta}{9c \frac{3+\delta/2}{1-\delta/2} \log(2/\delta)}. \quad (11)$$

Note that this condition on r/m is obtained from the extremal case $\ell = 1$, which minimizes the term $\left(\frac{1+\delta}{1-(1-\varepsilon)2^{-\ell-1}} - 16c \frac{\log(1/\varepsilon)}{1-\varepsilon} \ell^4 \cdot \left(\frac{2r}{m}\right)^\ell \right)$ in (10). Thus,

$$(8) = -\Omega \left((1-\varepsilon)2^{-\ell-1} \binom{m}{\leq r} \right) \leq -\Omega \left(2^{-r} \binom{m}{\leq r} \right).$$

Concluding, the overall probability $\Pr\{\mathcal{B}\}$ is bounded as

$$\Pr\{\mathcal{B}\} \leq 2^{-\Omega\left(\binom{m}{\leq r} 2^{-r}\right)} + 2^{-\Omega\left(\delta\binom{m}{\leq r}\right)} = 2^{-\Omega\left(\min(\delta, 2^{-r})\binom{m}{\leq r}\right)}.$$

□

6.2 High-rate regime

In this section we prove our main result for the BSC in the high rate regime.

Theorem 6.2. *Let $\varepsilon > 0$, $r \leq m$ two positive integers and $s = \lfloor \binom{m-\log\left(\binom{m}{\leq r}\right)+\log(\varepsilon)}{r} \rfloor - 1$. Then $RM(m, m - (2r + 2))$ can correct a random error pattern of weight s with probability larger than $1 - \varepsilon$.*

Using Claim 5.3, Theorem 6.2 gives the following corollary.

Corollary 6.3. *Let $\varepsilon, \delta > 0$, $r \leq m$ two positive integers such that $r < \sqrt{\frac{\delta m}{4 \log(m)}}$ and $\varepsilon > m^{-r/2}$. Then $RM(m, m - (2r + 2))$ can correct a random error pattern of weight $(1 - \delta) \binom{m}{\leq r}$ with probability larger than $1 - \varepsilon$.*

The rest of this section is organized as follows. We first give a combinatorial view of the syndrome of an error pattern under $E(m, r)$ (Section 6.2.1). We then study the case of $E(m, 3)$, which corresponds to the case $r = 1$ in Theorem 6.2 (as $E(m, 3) = H(m, m - 4)$), in Section 6.2.2. The case of general degree- r is handled in Section 6.2.3. Then, in Section 6.2.4 we extend the case $r = 1$ to hold for arbitrary linear codes of high degree and in Section 6.2.5 we prove that our results for the case $r = 1$ are tight, in some sense.

6.2.1 Parity check matrix and parity of patterns

In this section we give a combinatorial interpretation of the syndrome of an error pattern. Consider the code $RM(m, m - r - 1)$. Its parity check matrix is $H(m, m - r - 1) = E(m, r)$.

Let $U \subseteq \mathbb{F}_2^m$ be a set of size s . We associate with U the error pattern $\mathbb{1}_U \in \mathbb{F}_2^n$. Clearly $w(\mathbb{1}_U) = |U| = s$. We denote with u_j the j 'th element of U . We shall also think of U as an $m \times s$ matrix whose j 'th column is u_j . As before we denote with U^r the submatrix of $E(m, r)$ whose columns are indexed by U . Alternatively, this is the set of all evaluation vectors of U 's columns. We shall use the same convention for another subset $V \subseteq \mathbb{F}_2^m$.

The following definition captures a combinatorial property that we will later show its relation to syndromes under $E(m, r)$.

Definition 6.4. *For two matrices A, B of same dimension $n_1 \times n_2$, we denote $A \sim_r B$ if any pattern of size at most r in the columns of A appears with the same parity in the columns of B . I.e., for every subset $I \subset [n_1]$ of size r and every $z \in \mathbb{F}_2^r$ the number of columns in A_I , that equal z is equal, modulo 2, to the number of columns in B_I , that equal z .*

For example, the matrices

$$A = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \quad \text{and} \quad B = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 & 1 \end{pmatrix} \quad (12)$$

satisfy $A \sim_2 B$ but $A \not\sim_3 B$. To see that $A \sim_2 B$ one can observe that: the number of 1's in row i in A is equal (modulo 2) to that number in B ; the inner product (modulo 2) between rows i and j in A is the same as in B . Indeed, that inner product between rows i and j counts the number of columns that have 1 in both rows. Together with the information about the number of 1's in row i and in row j we are guaranteed that any pattern on rows i and j has the same parity in both matrices. On the other hand, the pattern $(1, 1, *, *, 1, *)$, which stands for 1 in the first, second and fifth rows (in the terminology of the definition, $I = \{1, 2, 5\}$ and $z = (1, 1, 1)$), appears once in B but it does not appear in A .

The next lemma shows that two error patterns $\mathbb{1}_V$ and $\mathbb{1}_U$ have the same syndrome under $E(m, r)$ if and only if the two matrices U and V satisfy $U \sim_r V$. We denote with $\mathbb{M}(m, r)$ the set of all m -variate monomials of degree at most r .

Lemma 6.5 (Parity of patterns). *For two sets $U, V \subseteq \mathbb{F}_2^m$ of size s it hold that*

$$E(m, r) \cdot \mathbb{1}_U = E(m, r) \cdot \mathbb{1}_V \iff \sum_{i=1}^s f(u_i) = \sum_{i=1}^s f(v_i), \quad \forall f \in \mathbb{M}(m, r) \quad (13)$$

$$\iff \sum_{i=1}^s f(u_i) = \sum_{i=1}^s f(v_i), \quad \forall f \in \mathbb{P}(m, r), \quad (14)$$

$$\iff U \sim_r V. \quad (15)$$

Proof. The first equivalence is by definition and the second one is clear. The equivalence between (14) and (15) is best explained by the following example. Consider the polynomial $f(x_1, \dots, x_m) = x_1(1 + x_m)$. In order for the equivalence to hold, it must be that the number of u_i 's (which are m -bit vectors themselves) that have 1 in the first coordinate and 0 in the last coordinate, is equal (modulo 2) to the number of v_i 's that have the same structure. In other word, the equation $\sum_{i=1}^s f(u_i) = \sum_{i=1}^s f(v_i)$ makes sure that the number of columns that have 1 in the first row and 0 in the last row is the same (modulo 2) in U and in V . By choosing different polynomials of degree at most r , the same must hold for any pattern of size at most r , and hence, by definition, $U \sim_r V$. The reverse direction is proved in a similar way. We note that the formal proof follows by induction on r . We leave the exact details of the proof to the reader. \square

Our goal is to understand, for given values of m and r , how many vectors $\mathbb{1}_U$ are bad, in the sense that they admit a bad companion $\mathbb{1}_V$ such that $E(m, r) \cdot \mathbb{1}_U = E(m, r) \cdot \mathbb{1}_V$. Thus, by Lemma 6.5 this is equivalent to studying pairs U, V such that $U \sim_r V$. The next lemma will allow us to apply linear transformations to U in order to make it “nicer” without losing generality.

Lemma 6.6 (Affine invariance). *In the notation of Lemma 6.5, if $U \sim_r V$ then $(AU) \sim_r (AV)$ for any linear transformation $A : \mathbb{F}_2^m \rightarrow \mathbb{F}_2^m$. Furthermore, if A is invertible then $(AU) \sim_r (AV)$ implies that $U \sim_r V$.*

Proof. By Lemma 6.5, $AU \sim_r AV$ iff $\sum_{i=1}^s f(u_i) = \sum_{i=1}^s f(v_i)$, $\forall f \in \mathbb{P}(m, r)$. For a polynomial f denote $f_A(x) \triangleq f(Ax)$. It is clear that $\deg(f_A) \leq \deg(f)$. It thus follows that

$$\begin{aligned} U \sim_r V &\iff \sum_{i=1}^s f(u_i) = \sum_{i=1}^s f(v_i) \quad \forall f \in \mathbb{P}(m, r), \\ &\implies \sum_{i=1}^s f_A(u_i) = \sum_{i=1}^s f_A(v_i), \quad \forall f \in \mathbb{P}(m, r), \\ &\iff \sum_{i=1}^s f(Au_i) = \sum_{i=1}^s f(Av_i) \quad \forall f \in \mathbb{P}(m, r), \\ &\iff AU \sim_r AV. \end{aligned}$$

To see the furthermore part, we note that if A is invertible then $f_{A^{-1}}(AV) = f(V)$. Hence, all the implications above can be made “if and only if”. \square

6.2.2 The case $r = 1$

To prove Theorem 6.2 we first study the case where $r = 1$ as the proof for the general case will use ideas similar to the proof of this case. Note that this case corresponds to studying syndrome

of error patterns under $H(m, m - (2r + 2)) = H(m, m - 4) = E(m, 3)$, namely evaluations by all degree-3 monomials.

We will prove first a deterministic result: if U is *any* set of linearly independent columns, then for any $V \neq U$, we have that $U \not\sim_3 V$. Thus, any set of errors that is supported on linearly independent coordinates (when viewed as vectors in \mathbb{F}_2^m) can be uniquely corrected. This immediately gives an average-case result. If we have $m - \log(m/\varepsilon)$ random errors, then with probability at least $1 - \varepsilon$ their locations correspond to linearly independent m -bit vectors and therefore we can correct such amount of errors with high probability.²⁰ Notice that this is already highly nontrivial, as $R(m, m - 4)$ has (absolute) distance 16, so in the worst case one cannot correct more than 8 worst-case errors!

Lemma 6.7. *Let $U \subseteq \mathbb{F}_2^m$ be a set of linearly independent vectors, such that $|U| = s$. Then, for any $V \neq U$, such that $|V| \leq s$, we have that $V \not\sim_3 U$.*

In particular this means that we can correct the error pattern $\mathbb{1}_U$ in $RM(m, m - 4)$.

Proof. By multiplying U with an invertible A (changing the basis \mathbb{F}_2^m) we can assume, w.l.o.g., that the columns of U are the elementary basis vectors, e_1, \dots, e_s .²¹ Indeed, since A is invertible it follows from Lemma 6.6 that it is enough to prove the claim for AU .

Let $V \subseteq \mathbb{F}_2^m$ be such that $|V| = s$ and $V \sim_3 U$. Our task is to show that $V = U$. This will be shown in two steps. First, we'll show that $\text{span}(V) = \text{span}(U)$, which in particular implies that V is linearly independent as well. Proving linear independence requires only that $V \sim_2 U$, namely evaluations by degree-2 monomials. Using $V \sim_3 U$, we'll prove that they actually have the same span, and from that derive that $V = U$.

Let us first argue linear independence of V . We'll think of U and V not only as sets of vectors, but also as $m \times s$ matrices, and denote by U' the transpose of U . Note that, as the columns of U are unit vectors, we have $U'U = I_s$. Now since diagonal elements of this product capture the value of degree-1 monomials of the syndrome, and off-diagonal elements of the of the product correspond to inner products of rows, namely (as in the example of the previous section), to degree-2 monomials of the syndrome. As $V \sim_2 U$ we also have that $V'V = I_s$ and so the dimension of V is s as well. We will later show 6.2.5 that this linear independence is the only thing we can infer from $V \sim_2 U$.

We now actually prove the stronger statement that in fact U and V span the same subspace. This will require $V \sim_3 U$. It will be sufficient to prove that V spans the vector e_1 , as for other vectors in U the proof is identical. Consider the pattern $(1, 0)$ in the first two rows of U . That is, consider all columns of U that have 1 in their first coordinate and 0 in the second. It is clear that this pattern only appears in e_1 and hence its parity in U is 1. Thus, there must be an odd number of columns in V whose first two rows equal $(1, 0)$. The main observation is that if we add up the columns then we obtain the vector e_1 .

Claim 6.8. *Under the conditions of the lemma, the sum of all columns in V whose first two coordinates equal $(1, 0)$ is e_1 . More generally, for $i \in [s]$, if we consider the pattern that has 1 in the i 'th coordinate and 0 in some $j \neq i$ coordinate, then the sum of all columns in V that have this pattern is equal to e_i .*

Proof. Assume that this is not the case, namely, the sum is a vector $w \neq e_1$. We first note that the first two coordinates of w equal $(1, 0)$. Indeed, this holds as we summed an odd number of vectors that has these values. Hence, there must exist a coordinate $i > 2$ such that $w_i = 1$. Thus, the

²⁰To eliminate possible confusion we repeat: an error pattern is an n -bit vector, whose coordinates are indexed by m -bit vectors.

²¹This is not really necessary, but it makes the argument simpler to explain.

number of vectors in V with the pattern $(1,0,1)$ in rows $(1,2,i)$ is odd. But this is not the case in U , contradicting $V \sim_3 U$.

The proof of the general case is similar. \square

We now use the fact that U and V have the same span to conclude that $U = V$. Denote $J_{10} \subset [s]$ the indices of columns in V that have $(1,0)$ as their first two coordinates. By Claim 6.8 we have that $\sum_{i \in J_{10}} v_i = e_1$. Next, consider the pattern $(1,*,0)$, namely, the pattern that has 1 in the first row and 0 in the third row. Denote the corresponding set of column indices with J_{1*0} . Again, Claim 6.8 implies that $\sum_{i \in J_{1*0}} v_i = e_1$. However, since the columns in V are linearly independent, there is only one way to represent e_1 as a linear combination of the columns of V and therefore it must be the case that $J_{10} = J_{1*0}$.

Continuing in this fashion we get that $J_{10} = J_{1*0} = J_{1**0} = \dots = J_{1***0}$, where the last set corresponds to the columns that have 1 in the first coordinate and 0 in the last coordinate. As the size of J_{10} is odd we know that it is not empty. In particular, all the vectors in J_{10} must satisfy that they have 1 in the first coordinate and 0 in the remaining coordinates. Indeed each such 0 can be justified by one of those J sets. In particular e_1 is a column in V . Repeating this process for all e_i , $i \in [s]$, we get that all these e_i 's are columns in V . Since V has exactly s columns it must have the same set of columns as U . In particular, $V = U$. This completes the proof of Lemma 6.7. \square

Lemma 6.7 shows that if the columns of U are linearly independent then $E(m,3)\mathbb{1}_U \neq E(m,3)\mathbb{1}_V$ for any other V of the same size. As randomly picking $m - \log(m/\varepsilon)$ vectors in \mathbb{F}_2^m we are likely to get linearly independent vectors we obtain our main result for the case $r = 1$.

Claim 6.9. *Let $\varepsilon > 0$ and $t = m - \log(m/\varepsilon)$. Pick t vectors uniformly at random from \mathbb{F}_2^m , $u_1, \dots, u_t \in \mathbb{F}_2^m$. Then, with probability at least $1 - \varepsilon$, the u_i are linearly independent.*

Proof. Set $u_0 = \vec{0}$. For each $1 \leq i \leq t$, the probability that u_i belongs to the span of u_0, \dots, u_{i-1} is at most $2^i / 2^m \leq \frac{\varepsilon}{m}$. The claim now follows from the union bound. \square

Combining Lemma 6.7 with Claim 6.9 we obtain the following corollary which is a special case of our main theorem.

Corollary 6.10. *We can correct a random set of $m - \log(m/\varepsilon)$ errors in $RM(m, m - 4)$ with high probability.*

6.2.3 The degree- r case

The proof of the degree- r case proceeds along the same lines as the proof of the $r = 1$ case. However, in order to correct errors for $RM(m, m - (2r + 2))$ we will require that the matrix U^r corresponding to the error pattern $\mathbb{1}_U$ has linearly independent columns. Note that when $r = 1$ this amounts to requiring that U has linearly independent columns, in order to correct $\mathbb{1}_U$ in $RM(m, m - 4)$, which is exactly what we proved in Section 6.2.2. This is also a shortcoming of our result, it is clear that with this condition we cannot expect to correct more than $\binom{m}{\leq r}$ errors. On the one hand this is much better than the minimum-distance based result of 2^{2r+2} , but on the other hand one may hope to be able to decode from $O\left(\binom{m}{\leq 2r+1}\right)$ many errors.

Our main lemma is an analog of Lemma 6.7.

Lemma 6.11. *Let $U \subseteq \mathbb{F}_2^m$ be such that $|U| = s$ and the columns of U^r are linearly independent. Then, for any $V \subseteq \mathbb{F}_2^m$ satisfying $|V| = s' \leq s$ and $V \neq U$, we have that $V \not\sim_{2r+1} U$.*

Notice that Lemma 6.7 is obtained by setting $r = 1$ in Lemma 6.11.

Proof. Denote with u_i the i 'th column of U . Recall that the i 'th column of U^r corresponds to all evaluations of monomials of degree at most r at the point u_i , i.e., it is equal to u_i^r . This interpretation will be helpful throughout the proof. Assume that $V \subseteq \mathbb{F}_2^m$ is such that $|V| = s' \leq s$ and $V \sim_{2r+1} U$. Similarly, we denote the columns of V with $v_1, \dots, v_{s'} \in \mathbb{F}_2^m$ and note that the i 'th column of V^r is v_i^r .

As the columns of U^r are linearly independent, there exist vectors f_i so that $f_i \cdot U^r = e_i \in \mathbb{F}_2^s$. As the coordinates of each f_i are indexed by monomials of degree $\leq r$, we can interpret f_i as a degree- r polynomial $f_i(x_1, \dots, x_m)$ and rewrite $f_i \cdot U^r$ as $f_i \cdot U^r = (f_i(u_1), \dots, f_i(u_s))$. In other words, $f_i(u_j) = \delta_{i,j}$.²²

Our next goal is proving that if $U \sim_{2r+1} V$ then $u_1 \in V$. This will clearly imply the lemma as we can prove the same for any other u_i . Our main handle will be the polynomial f_1 that separates u_1 from the other u_i 's.

Let us assume wlog that $(u_1)_1 = 1$, i.e., that the first coordinate of u_1 equals 1.²³ Consider the polynomial $x_1 \cdot f_1$. This is a polynomial of degree $\leq r+1$ and by the definition of f_1 and the assumption on $(u_1)_1$ we have that $\sum_{i=1}^s (x_1 f_1)(u_i) = 1$. As $U \sim_{2r+1} V$, it must hold that $\sum_{i=1}^{s'} (x_1 f_1)(v_i) = 1 \pmod{2}$.

Following the footsteps of the proof of Lemma 6.7, we denote $J_{f_1,1} = \{i \mid (x_1 f_1)(v_i) = 1\}$. The next claim is analogous to Claim 6.8.

Claim 6.12.

$$\sum_{i \in J_{f_1,1}} v_i^r = u_1^r.$$

Proof. Let M be some monomial of degree $\leq r$. To ease the notation assume that $M(u_1) = 1$ and consider the polynomial $M \cdot x_1 \cdot f_1$.²⁴ It is clear that $(M \cdot x_1 \cdot f_1)(u_1) = 1$. Since $V \sim_{2r+1} U$ and $\deg(M \cdot x_1 \cdot f_1) \leq 2r+1$, it follows that $\sum_{i=1}^{s'} (M x_1 f_1)(v_i) = 1$. From definition of $J_{f_1,1}$ we have that

$$\sum_{i=1}^{s'} (M x_1 f_1)(v_i) = \sum_{i \in J_{f_1,1}} (M x_1 f_1)(v_i) = 1.$$

Indeed, for every $i \notin J_{f_1,1}$ we have that $(x_1 f_1)(v_i) = 0$. We thus conclude that there is an odd number of vectors v_i , $i \in J_{f_1,1}$, such that $(M x_1 f_1)(v_i) = 1$. In particular, the M 'th coordinate in the sum $\sum_{i \in J_{f_1,1}} v_i^r$ equals 1, i.e. it is equal to $M(u_1)$. As M was arbitrary we conclude that $\sum_{i \in J_{f_1,1}} v_i^r = u_1^r$, as required. \square

As we can prove an analogous lemma for every u_i , we conclude the the columns of U^r belong to the span of the columns of V^r . In particular, the columns of V^r are linearly independent and $|V| = s$ (earlier we called this observation Claim ??). Our next step is proving that, up to permutation of columns, $U^r = V^r$. This will imply that $u_i = v_i$ as we wanted.

To show this, for every $\ell \in [m]$, we denote

$$J_{f_1,\ell} = \{i \mid (v_i)_\ell = (u_1)_\ell \text{ and } f_1(v_i) = 1\}.$$

²²Intuitively, the polynomials f_i correspond to the rows of the matrix A that were used in the proof of Lemma 6.7 to make the columns of U equal the elementary unit vectors there.

²³If it equals 0 then we consider the polynomial $(1 + x_1)f_1$ in what follows.

²⁴If $M(u_1) = 0$ then we consider the polynomial $(1 + M) \cdot x_1 \cdot f_1$ instead.

We note that there is an alternative way to define $J_{f_1, \ell}$ by considering either the polynomial $x_\ell \cdot f_1$ or the polynomial $(1 + x_\ell) \cdot f_1$. We thus have that for every $\ell \in [m]$,

$$\sum_{i \in J_{f_1, \ell}} v_i^r = u_1^r.$$

However, as the columns of V^r are linearly independent and

$$\sum_{i \in J_{f_1, \ell}} v_i^r = u_1^r = \sum_{i \in J_{f_1, 1}} v_i^r$$

we get that $J_{f_1, 1} = J_{f_1, 2} = \dots = J_{f_1, m}$. Hence,

$$J_{f_1, 1} = \cap_{\ell=1}^m J_{f_1, \ell} = \{i \mid \forall \ell \in [m] (v_i)_\ell = (u_1)_\ell \text{ and } f_1(v_i) = 1\}.$$

Thus, for every $i \in J_{f_1, 1}$ we have that $v_i = u_1$. In particular, since $J_{f_1, 1} \neq \emptyset$, it follows that there is some $i \in [s]$ such that $v_i = u_1$. As we can prove this for every u_j , we conclude that $U = V$ as claimed. This concludes the proof of Lemma 6.11. \square

We thus proved that if an error pattern $\mathbb{1}_U$ is such that its coordinates u_i satisfy that u_i^r are linearly independent, then we can correct that error pattern in $RM(m, m - (2r + 2))$. We summarise this in the following theorem.

Theorem 6.13. *If a set of columns U are linearly independent in $E(m, r)$ (namely, $RM(m, m - r - 1)$ can correct the erasure pattern $\mathbb{1}_U$), then the error pattern $\mathbb{1}_U$ can be corrected in $RM(m, m - (2r + 2))$.*

Proof. Lemma 6.11 tells us that if the columns indexed by U are linearly independent in $E(m, r)$, then there is no other $V \subseteq \mathbb{F}_2^m$ of size $\leq s$ such that $V \sim_{2r+1} U$. Lemma 6.5 now implies that

$$E(m, 2r + 1) \cdot \mathbb{1}_U \neq E(m, 2r + 1) \cdot \mathbb{1}_V,$$

for any $U \neq V \subseteq \mathbb{F}_2^m$ of size $\leq s$. As $E(m, 2r + 1) = H(m, m - (2r + 2))$, it follows that the syndrome of $\mathbb{1}_U$ is unique and hence $\mathbb{1}_U$ is uniquely decodable in $RM(m, m - (2r + 2))$. \square

The proof of Theorem 6.2 immediately follows from Theorems 4.5 and 6.13.

Proof of Theorem 6.2. Theorem 4.5 guarantees that a set $U \subseteq \mathbb{F}_2^m$ of $s = \lfloor \binom{m - \log(\frac{m}{\leq r}) + \log(\varepsilon)}{\leq r} \rfloor - 1$ randomly chosen vectors, satisfy that the columns of U^r are linearly independent. By Lemma 6.11 we learn that there is not other $V \subseteq \mathbb{F}_2^m$ of size $\leq s$ such that $V \sim_{2r+1} U$. Lemma 6.5 implies that for any such V ,

$$E(m, 2r + 1) \cdot \mathbb{1}_U \neq E(m, 2r + 1) \cdot \mathbb{1}_V.$$

As $E(m, 2r + 1) = H(m, m - (2r + 2))$, it follows that the syndrome of $\mathbb{1}_U$ is unique and hence $\mathbb{1}_U$ is uniquely decodable in $RM(m, m - (2r + 2))$. \square

6.2.4 A general reduction from decoding from errors to decoding from erasures

In this section we show that the results proved in Section 6.2.2 are in fact more general and apply to any degree three tensoring of a linear code with itself. We first set up the required definitions.

Definition 6.14. *The Hadamard product of two vectors $y, z \in \mathbb{F}_2^n$ is the vector $w = y \circ z$ obtained from the coordinate wise product $w_i = y_i \cdot z_i$.*

Definition 6.15. Let H be a $k \times n$ matrix. For every natural number ℓ , $H^{\otimes \ell}$ is a $\binom{k}{\leq \ell} \times n$ matrix that is defined as follows. Rows of $H^{\otimes \ell}$ are indexed by tuples $i_1 < i_2 < \dots < i_j$, $1 \leq j \leq \ell$, where the corresponding row in $H^{\otimes \ell}$ is equal to the Hadamard product of rows i_1, i_2, \dots, i_j .

In other words, if we think of H as the set of its column vectors then, using our usual notation, $H^{\otimes \ell} = H^\ell$. In particular, the parity check matrix $H(m, m - r - 1)$ of the code $RM(m, m - r - 1)$ is equal to $H(m, 1)^r$. Indeed the row indexed by $i_1 < i_2 < \dots < i_j$ corresponds to the evaluations of the monomial $\prod_{t=1}^j x_{i_t}$.

It is also clear that for any integers m, n and any $m \times n$ matrix H , the set of columns of H is contained in the set of columns of the Hadamard matrix of rank m , i.e., $E(m, 1)$, namely, every column of H appears in $E(m, 1)$. We thus obtain the following two corollaries that follow from the proof technique of the previous section.²⁵

Corollary 6.16. Let m, n be integers and H an $m \times n$ matrix. Let $S \subseteq [n]$ be such that the columns indexed by S in H are linearly independent. Then, in the code whose parity check matrix is $H^{\otimes 3}$, we can correct the error pattern S .

Using the relation between correcting erasures and independence (Lemma 2.8) in the parity check matrix we obtain the following corollary.

Theorem 6.17. Let $C \subseteq \mathbb{F}_2^n$ be a linear code with parity check matrix H . For any subset $S \subseteq [n]$ the following hold: If we can recover codewords in C from erasures in the coordinates S then in the code whose parity check matrix is $H^{\otimes 3}$, we can correct the error pattern $\mathbb{1}_S$.

We note that Corollary 6.10 is a special case of Theorem 6.17.

6.2.5 The degree-2 counterexample

In this section we prove that $RM(m, m - 3)$ does not achieve capacity for the BSC. For this Reed-Muller code to achieve capacity we must have that, w.h.p., a random error pattern of weight $O(m)$ has a unique syndrome. We next show that very few patterns of this weight have unique syndromes. In fact, we show that very few patterns of weight \sqrt{m} have unique syndromes.

Let $s < \sqrt{m}$ be an even integer. Let B be the following $s \times s$ matrix. For $1 \leq i \leq s - 2$, the i 'th row of B has 1 in coordinates 1, 2 and $i + 2$. E.g., the first row of B begins with three 1's followed by zeros. The $s - 1$ 'th row of B equals $(1, 0, 1, \dots, 1)$ and the last row of B is $(0, 1, \dots, 1)$. Note that $B \cdot B^t = I$. For example, the matrix B in (12) is what we get if we set $s = 6$.

Let u_1, \dots, u_s be any set of vectors in \mathbb{F}_2^m . Let U be the matrix whose i 'th column is u_i . Define $V = U \cdot B$. It is not hard to verify that $E(m, 2) \cdot \mathbb{1}_U = E(m, 2) \cdot \mathbb{1}_V$. Thus, if $V \neq U$ then $\mathbb{1}_U$ does not have a unique syndrome.

Finally, we note that picking u_1, \dots, u_s at random is equivalent to picking the matrix U at random. If U and V have the same set of columns then there must exist an $s \times s$ permutation matrix Π such that $U\Pi = V$. Thus, $U(B - \Pi) = 0$. Fix Π . The probability that all rows of U are in the kernel of $B - \Pi$ is at most 2^{-m} . Indeed, for every permutation matrix Π , $\text{rank}(B - \Pi) \geq 1$. As there are $s!$ permutation matrices, the probability that U is unique is at most $s!/2^m < 2^{-m/2}$.

²⁵The proofs are completely identical and are thus omitted.

7 Future directions and open problems

We believe that our work renews hope for progress on some classical questions, and suggests some new concrete directions and open problems.

The most obvious of all is the question of whether Reed-Muller codes achieve capacity for all ranges of parameters, either for random erasures or for random errors. We only handle here the extreme cases of very high or very low rates, whereas most interest is traditionally focused on constant rate codes. We believe that the techniques for each of our four bounds can be improved to a larger set of parameters (see below), but feel that they fall short of reaching constant rate, and possibly new techniques are needed.

One way to improve our bounds in both Theorem 1.1 (low-rate BEC) and Theorem 1.7 (low rate BSC) is through tighter bounds on the weight enumeration of Reed-Muller codes, as well as tighter bounds on the probability of error for Theorem 1.7. We believe that in Theorem 1.5 one can eliminate the factor ℓ^4 in the exponent, resulting in a bound that is a fixed polynomial (independent of m, r, ℓ) of the lower bound in [KLP12]. While such a tight result would not get us (in either Theorem 1.1 and 1.7) to the constant rate regime, this question of weight enumeration is of course basic in its own right. Moreover, both in [KLP12] and our paper, it also implies similar bounds for list-decoding, which is another basic question.

Theorem 1.4 (high rate BEC) is quantitatively much weaker than Theorems 1.1 and 1.7, in that the latter two can handle polynomials of degree- r which is linear in m , whereas the former only reaches degrees r which are about \sqrt{m} . The bottleneck in the argument, which probably prevents it from reaching a linear degree, is the use of the union bound. We upper bound the probability that, when adding a subsequent random vector u to our set U , its evaluation u^r will be linearly independent of the evaluations of all previously chosen points. This current proof does not use at all that previous points were chosen randomly, as we don't know how to take advantage of this.

For high-rate BSC (Theorem 1.9), while we are able to correct many more errors than previously known, we are not even able to achieve capacity. Here we feel that one important bottleneck is our inability to argue directly about corruption patterns (sets U) which are linearly *dependent*. Our unique decoding proof, even for $r = 1$ (on which we focus now), showing that a set $U \in \mathbb{F}_2^m$ is uniquely determined by its syndrome under evaluations by degree-3 monomials i.e., by $E(m, 3) \cdot \mathbf{1}_U$, is especially tailored to linearly independent sets U . The gap between our lower bound (namely that $E(m, 2) \cdot \mathbf{1}_U$ does not suffice) and the above upper bound (that $E(m, 3) \cdot \mathbf{1}_U$ suffices) is intriguing, and we believe we can find a subset of quadratically many monomials of degree at most 3 which guarantee unique decoding - such a result is information theoretically optimal; number of error patterns U which are linearly independent is about $\exp(m^2)$, and thus $O(m^2)$ bits are needed in any unique encoding.

Another burning question regarding this result is its inefficiency. While unique decoding is guaranteed, the best way we know to identify the set U is brute force, requiring $\exp(m^2)$ steps for independent sets U of size m . We feel that a good starting place is (perhaps using our uniqueness proof) which recovers U in $\exp(m) = \text{poly}(n)$ steps from its evaluation on all degree-3 monomials (or even degree-10 monomials). Of course, it is quite possible that a $\text{poly}(m)$ algorithm exists. In particular, recursive algorithms (that exploit the recursive nature of RM codes) could be used to that effect²⁶.

²⁶Practically, one can decode RM codes on the BSC by using a recursive decoder (e.g., like for polar codes) for each missing component in the syndrome, and by growing a list of possible codewords each time the decoder has doubts, or pruning down the tree each time the decoder can check the validity of a path (from the available components of the syndrome).

Yet another research direction related to our result is of course exploring the connections between recovering from erasures and from errors. Our general reduction between the two uses tensor powers and hence loses in efficiency (which here is best captured by the co-dimension of the code, which is cubed). Is there a reduction which loses less? We do not know how to rule out a reduction that increases the co-dimension only by a constant factor. There is no reason to restrict attention to Reed-Muller codes and our tensor construction - such a result would be of use anywhere, as erasures are so much simpler to handle than errors.

Finally, we believe that a better understanding of the relation between Reed-Muller codes and Polar codes is needed, and perhaps more generally an understanding of which subspaces of polynomials generated by subsets of monomials give rise to good, efficient codes. In particular, it would also be interesting to investigate the scaling of the blocklength in terms of the gap to capacity for RM codes. It was proved recently in [GX13] that for polar codes, the blocklength scales polynomially with the inverse of the gap to capacity, with a precise characterization given in [Has13]. While this scaling does not match the optimal scaling of random codes [Str62], it is in contrast to the exponential scaling obtained with concatenated codes [For67] (see [GX13] for a discussion on this). It would be interesting to investigate such finer questions for RM codes in view of the results obtained in this paper, which already provide partial information about these scalings.

Acknowledgements

We thank Venkatesan Guruswami for bringing [Wei91] to our attention. The second author would like to thank the organizers of Dagstuhl meeting “Algebra in Computational Complexity,” where he discussed the results with Venkatesan Guruswami.

References

- [Abb11] E. Abbe, *Randomness and dependencies extraction via polarization*, ITA, 2011, Available at arXiv:1102.1247, pp. 1–7. [3](#)
- [AL04] A. Ashikhmin and S. Litsyn, *Simple MAP decoding of first-order Reed-Muller and Hamming codes*, Information Theory, IEEE Transactions on **50** (2004), no. 8, 1812–1818. [11](#)
- [ALM⁺98] S. Arora, C. Lund, R. Motwani, M. Sudan, and M. Szegedy, *Proof verification and the hardness of approximation problems*, Journal of the ACM (JACM) **45** (1998), no. 3, 501–555. [1](#)
- [Ari08] E. Arikan, *A performance comparison of polar codes and Reed-Muller codes*, Communications Letters, IEEE **12** (2008), no. 6, 447–449. [2](#), [3](#)
- [Ari09] ———, *Channel polarization: A method for constructing capacity-achieving codes for symmetric binary-input memoryless channels*, Information Theory, IEEE Transactions on **55** (2009), no. 7, 3051–3073. [2](#), [3](#), [12](#)
- [AS92] N. Alon and J. Spencer, *The probabilistic method*, John Wiley, 1992. [5](#)
- [BCM⁺V14] A. Bhaskara, M. Charikar, A. Moitra, and A. Vijayaraghavan, *Smoothed analysis of tensor decompositions*, Proceedings of the 46th Annual ACM Symposium on Theory of Computing, STOC ’14, 2014. [21](#)

- [BF90] D. Beaver and J. Feigenbaum, *Hiding instances in multioracle queries*, STACS 90, Springer, 1990, pp. 37–48. [1](#)
- [BFL90] L. Babai, L. Fortnow, and C. Lund, *Nondeterministic exponential time has two-prover interactive protocols*, Foundations of Computer Science, 1990. Proceedings., 31st Annual Symposium on, IEEE, 1990, pp. 16–25. [1](#)
- [BGH⁺12] B. Barak, P. Gopalan, J. Håstad, R. Meka, P. Raghavendra, and D. Steurer, *Making the long code shorter*, 53rd Annual IEEE Symposium on Foundations of Computer Science, FOCS 2012, New Brunswick, NJ, USA, October 20–23, 2012, 2012, pp. 370–379. [1](#)
- [BKS⁺10] A. Bhattacharyya, S. Kopparty, G. Schoenebeck, M. Sudan, and D. Zuckerman, *Optimal testing of reed-muller codes*, Proceedings of the 2010 IEEE 51st Annual Symposium on Foundations of Computer Science, FOCS '10, 2010, pp. 488–497. [1](#)
- [BV10] A. Bogdanov and E. Viola, *Pseudorandom bits for polynomials*, SIAM J. Comput. **39** (2010), no. 6, 2464–2486. [1](#)
- [DG07] D.J. Costello, Jr. and G. D. Forney, Jr., *Channel coding: The road to channel capacity*, Proceedings of the IEEE **95** (2007), no. 6, 1150–1177. [2](#)
- [DS06] I. Dumer and K. Shabunov, *Recursive error correction for general Reed-Muller codes*, Discrete Applied Mathematics **154** (2006), no. 2, 253 – 269, Coding and Cryptography. [10](#)
- [Dum04] I. Dumer, *Recursive decoding and its performance for low-rate Reed-Muller codes*, Information Theory, IEEE Transactions on **50** (2004), no. 5, 811–823. [10](#)
- [Dum06] ———, *Soft-decision decoding of Reed-Muller codes: a simplified algorithm*, Information Theory, IEEE Transactions on **52** (2006), no. 3, 954–963. [10](#)
- [Eli55] P. Elias, *Coding for noisy channels*, IRE Convention Record **4** (1955), 37–46. [2](#)
- [For67] G. D. Forney, *Concatenated codes*, PhD Thesis, Massachusetts Institute of Technology, 1967. [2](#), [39](#)
- [Gas04] W. Gasarch, *A survey on private information retrieval*, Bulletin of the EATCS, Citeseer, 2004. [1](#)
- [GKZ08] P. Gopalan, A. R. Klivans, and D. Zuckerman, *List-decoding Reed-Muller codes over small fields*, Proceedings of the 40th Annual ACM Symposium on Theory of Computing, Victoria, British Columbia, Canada, May 17–20, 2008, 2008, pp. 265–274. [4](#)
- [GX13] V. Guruswami and P. Xia, *Polar codes: Speed of polarization and polynomial gap to capacity*, Proceedings of the 2013 IEEE 54th Annual Symposium on Foundations of Computer Science, 2013, pp. 310–319. [39](#)
- [Ham50] R. W. Hamming, *Error Detecting and Error Correcting Codes*, Bell System Technical Journal **26** (1950), no. 2, 147–160. [2](#)
- [Has13] S. H. Hassani, *Polarization and Spatial Coupling: Two Techniques to Boost Performance*, PhD Dissertation, IC, EPFL, 2013. [39](#)

- [HKL05] T. Helleseth, T. Klove, and V. I. Levenshtein, *Error-correction capability of binary linear codes*, Information Theory, IEEE Transactions on **51** (2005), no. 4, 1408–1423. [11](#)
- [ILL89] R. Impagliazzo, L. A. Levin, and M. Luby, *Pseudo-random generation from one-way functions (extended abstracts)*, Proceedings of the 21st Annual STOC, 1989, pp. 12–24. [43](#), [44](#)
- [KLP12] T. Kaufman, S. Lovett, and E. Porat, *Weight Distribution and List-Decoding Size of Reed-Muller Codes*, Information Theory, IEEE Transactions on **58** (2012), no. 5, 2689–2696. [1](#), [4](#), [7](#), [9](#), [10](#), [11](#), [17](#), [18](#), [19](#), [26](#), [38](#)
- [Kri70] R. E. Krichevskiy, *On the number of Reed-Muller code correctable errors*, Dokl. Sov. Acad. Sci. **191** (1970), 541–547. [10](#)
- [KRU11] S. Kudekar, T. J. Richardson, and R. L. Urbanke, *Threshold saturation via spatial coupling: Why convolutional LDPC ensembles perform so well over the BEC*, Information Theory, IEEE Transactions on **57** (2011), no. 2, 803–834. [2](#)
- [KT70] T. Kasami and N. Tokura, *On the weight structure of Reed-Muller codes*, Information Theory, IEEE Transactions on **16** (1970), no. 6, 752–759. [11](#)
- [KTA76] T. Kasami, N. Tokura, and S. Azumi, *On the weight enumeration of weights less than $2.5d$ of Reed-Muller codes*, Information and Control **30** (1976), no. 4, 380 – 395. [11](#)
- [LMS⁺97] M. Luby, M. Mitzenmacher, A. Shokrollahi, D. A. Spielman, and V. Stemann, *Practical loss-resilient codes*, 29th annual ACM Symposium on Theory of Computing, 1997, pp. 150–159. [2](#)
- [MHU14] M. Mondelli, S. H. Hassani, and R. L. Urbanke, *From Polar to Reed-Muller Codes: A Technique to Improve the Finite-Length Performance*, Communications, IEEE Transactions on **62** (2014), no. 9, 3084–3091. [2](#), [3](#)
- [MS77] F. J. MacWilliams and N. J. A. Sloane, *The theory of error correcting codes*, North-Holland mathematical library, no. v. 2, North-Holland Publishing Company, 1977. [11](#), [17](#)
- [Mul54] D. E. Muller, *Application of boolean algebra to switching circuit design and to error detection*, Electronic Computers, Transactions of the I.R.E. Professional Group on **EC-3** (1954), no. 3, 6–12. [1](#)
- [Rab89] M. O. Rabin, *Efficient dispersal of information for security, load balancing, and fault tolerance*, Journal of the ACM (JACM) **36** (1989), no. 2, 335–348. [1](#)
- [Raz87] A. A. Razborov, *Lower bounds on the size of bounded depth circuits over a complete basis with logical addition*, Math. Notes **41** (1987), no. 4, 333–338. [1](#)
- [Ree54] I. Reed, *A class of multiple-error-correcting codes and the decoding scheme*, Information Theory, Transactions of the IRE Professional Group on **4** (1954), no. 4, 38–49. [1](#), [10](#)
- [SB70] N. J. A. Sloane and E. Berlekamp, *Weight enumerator for second-order Reed-Muller codes*, Information Theory, IEEE Transactions on **16** (1970), no. 6, 745–751. [11](#)
- [Sha48] C. E. Shannon, *A mathematical theory of communication*, The Bell System Technical Journal **27** (1948), 379–423, 623–. [2](#)

- [Sha79] A. Shamir, *How to share a secret*, Communications of the ACM **22** (1979), no. 11, 612–613. [1](#)
- [Sha92] ———, *Ip= pspace*, Journal of the ACM (JACM) **39** (1992), no. 4, 869–877. [1](#)
- [Str62] V. Strassen, *Asymptotische Abschätzungen in Shannon's Informationstheories*, Trans. 3rd Prague Conf. Info. Theory (1962), 689–723. [39](#)
- [VMS92] A. S. Pershakov V. M. Sidel'nikov, *Decoding of Reed-Muller Codes with a Large Number of Errors*, Problems Inform. Transmission **28** (1992), no. 3, 80–94. [11](#)
- [Wei91] V. K. Wei, *Generalized hamming weights for linear codes*, IEEE Transactions on Information Theory **37** (1991), no. 5, 1412–1418. [9](#), [23](#), [24](#), [39](#)

A Proofs of Claim [4.15](#) and Claim [5.3](#)

Proof of Claim [4.15](#). We first need the following estimate.

Claim A.1. *For integers $0 \leq a \leq c$ and $b \leq c - a$ we have that $\sum_{i=1}^a \binom{c-i}{b} = \binom{c}{b+1} - \binom{c-a}{b+1}$.*

Proof.

$$\begin{aligned}
 \sum_{i=1}^a \binom{c-i}{b} &= \sum_{i=1}^a \binom{c-i}{b} + \binom{c-a}{b+1} - \binom{c-a}{b+1} \\
 &= \sum_{i=1}^{a-1} \binom{c-i}{b} + \binom{c-a+1}{b+1} - \binom{c-a}{b+1} \\
 &= \dots \\
 &= \binom{c}{b+1} - \binom{c-a}{b+1}.
 \end{aligned}$$

□

Using the claim we get that

$$\begin{aligned}
 \binom{m}{\leq r} - \sum_{i=1}^t \binom{m-i}{\leq r-1} &= \binom{m}{\leq r} - \sum_{i=1}^t \sum_{j=0}^{r-1} \binom{m-i}{j} \\
 &= \binom{m}{\leq r} - \sum_{j=0}^{r-1} \sum_{i=1}^t \binom{m-i}{j} \\
 &= \binom{m}{\leq r} - \sum_{j=0}^{r-1} \left(\binom{m}{j+1} - \binom{m-t}{j+1} \right) \\
 &= \binom{m-t}{\leq r}.
 \end{aligned}$$

□

Proof of Claim 5.3. We shall need the following two simple inequalities that hold for every $C > A/4 > B$:

$$\frac{\binom{A}{B}}{\binom{C}{B}} > \left(\frac{A-B}{C}\right)^B \quad \text{and} \quad \binom{A}{\leq B} \leq 2A^B.$$

Thus, for any $0 \leq r' \leq r$,

$$\begin{aligned} \frac{\binom{m-3\log(\binom{m}{\leq r})+\log(\varepsilon)}{r'}}{\binom{m}{r'}} &> \left(\frac{m-3\log(\binom{m}{\leq r})+\log(\varepsilon)-r'}{m}\right)^{r'} \\ &= \left(1 - \frac{3\log(\binom{m}{\leq r}) - \log(\varepsilon) + r'}{m}\right)^{r'} \\ &\geq \left(1 - \frac{3r\log(m) + 3 - \log(\varepsilon) + r'}{m}\right)^{r'} \\ &\geq \left(1 - \frac{4r\log(m)}{m}\right)^{r'} \\ &\geq \left(1 - \frac{4r' \cdot r\log(m)}{m}\right) \\ &\geq (1 - \delta). \end{aligned}$$

Hence,

$$\binom{m-3\log(\binom{m}{\leq r})+\log(\varepsilon)}{\leq r} = \sum_{r'=0}^r \binom{m-3\log(\binom{m}{\leq r})+\log(\varepsilon)}{r'} > \sum_{r'=0}^r (1 - \delta) \binom{m}{r'} = (1 - \delta) \binom{m}{\leq r},$$

as claimed. \square

B A proof of Lemma 4.10 using hashing

In this section we prove the following slightly weaker version of Lemma 4.10.

Lemma B.1. *Let $\mathcal{V} \subseteq \mathbb{F}_2^m$ such that $|\mathcal{V}| > 2^{m-t}$. Then there are more than $\binom{m-t-2\lceil\log(\binom{m-t}{\leq r})\rceil}{\leq r}$ linearly independent polynomials of degree $\leq r$ that are defined on \mathcal{V} .*

Notice that the only difference between Lemma 4.10 and Lemma B.1 is the constant 2 in the lower bound.

The main idea in the proof is showing that there exists a linear transformation T such that the projection of the set $T(\mathcal{V})$ onto (roughly) the first $\log(|\mathcal{V}|)$ coordinates contains a ball of radius r around some point. Since restricting monomials, of degree $\leq r$, to a ball of radius r yields linearly independent functions, the claim follows. To prove that a random transformation has a large projection onto the first coordinates we use the *leftover hash lemma* of Impagliazzo et al. [ILL89]. This is where we lose compared to Lemma 4.10. The lemma of [ILL89] gives more information than just a large projection (i.e., that the distribution on the projection is close to uniform) and so it does not get the same parameters that we can get using the result of Wei (Theorem 4.14).

Proof. We start by proving that, after a suitable linear transformation, the projection of \mathcal{V} onto the first coordinates contains a large ball.

Lemma B.2. *Let $Y \subseteq \mathbb{F}_2^m$ a set of size 2^a . Then, there is a linear transformation T such that for some $z \in \mathbb{F}_2^{a-2\lceil\log(\binom{a}{\leq r})\rceil}$, the ball $B(z, r)$ is contained in the projection of $T(A)$ onto the first $a - 2\lceil\log(\binom{a}{\leq r})\rceil$ coordinates.*

Proof. For the proof we need the following leftover hash lemma of Impagliazzo et al. [ILL89].

Lemma B.3 ([ILL89]). *Let $\ell \leq a \leq m$ be integers and $Y \subseteq \mathbb{F}_2^m$ a set of size 2^a . Then, there exists an invertible $m \times m$ matrix T such that the projection of the set $T(Y)$ onto the first $a - \ell$ coordinates yields a set of size larger than $2^{a-\ell}(1 - 2^{-\ell/2})$.*

Given Lemma B.3 the proof of Lemma B.2 is by a simple averaging argument. For the proof we shall denote with $\pi_b(\cdot)$ the map that projects m -bit vectors on their first b coordinates.

Denote $\hat{a} = a - 2\lceil\log(\binom{a}{\leq r})\rceil$. Apply Lemma B.3 with $\ell = a - \hat{a}$, a, m on the set Y . We get that, for a suitable linear transformation T , the projection $\pi_{\hat{a}}(T(Y))$ is a set of size larger than $2^{\hat{a}}(1 - 2^{-\lceil\log(\binom{a}{\leq r})\rceil}) \geq 2^{\hat{a}}(1 - \frac{1}{\binom{a}{\leq r}})$.

By linearity of expectation, there is a point $z \in \mathbb{F}_2^{\hat{a}}$ so that the fraction of points in $B(z, \hat{a}) \cap \pi_{\hat{a}}(T(Y))$ is larger than $|B(z, \hat{a})| \cdot (1 - \frac{1}{\binom{a}{\leq r}})$. As $|B(z, \hat{a})| = \binom{a}{\leq r}$, it follows that the size of the intersection is larger than $\binom{a}{\leq r}(1 - \frac{1}{\binom{a}{\leq r}}) = \binom{a}{\leq r} - 1$. Since the intersection size is an integer it must equal $\binom{a}{\leq r}$. In other words, $\pi_{\hat{a}}(T(Y))$ contains $B(z, r)$. \square

We continue with the proof of Lemma B.1. The point of the last two lemmas is that for such a set Y , the set of polynomials that are defined on it is isomorphic to the set of polynomials defined on $T(Y)$ (also when considering degree $\leq r$ polynomials for both sets). Let us focus on $T(Y)$ and consider only polynomials in the variables $x_1, \dots, x_{\hat{a}}$. As $\pi_{\hat{a}}(T(Y))$ contains a ball of radius r , we get that all monomials of degree $\leq r$ in $x_1, \dots, x_{\hat{a}}$ are linearly independent on $\pi_{\hat{a}}(T(Y))$. However, the value of any such monomial on a point in $T(Y)$ is the same as its value on its projection. Thus, there are at least $\binom{\hat{a}}{\leq r}$ many linearly independent polynomials, of degree $\leq r$, that are defined on $T(Y)$. Hence, there are at least $\binom{\hat{a}}{\leq r}$ many linearly independent polynomials, of degree $\leq r$, that are defined on Y .

In our case $|\mathcal{V}| > 2^{m-t}$. Thus, in the notation above, $a = m - t$ and $\hat{a} = a - 2\lceil\log(\binom{a}{\leq r})\rceil = m - t - 2\lceil\log(\binom{m-t}{\leq r})\rceil$. Thus, more than $\binom{m-t-2\lceil\log(\binom{m-t}{\leq r})\rceil}{\leq r}$ linearly independent polynomials of degree $\leq r$ that are defined on \mathcal{V} . \square